

2025 12 24
발표 자료

광운대학교 로봇학과
FAIR Lab

김한서

이번 주 진행사항

- ElasTST
 - 논문 리뷰
 - 실험 세팅
 - 실험 결과 및 시각화
 - 결과 정리

ElasTST: Towards Robust Varied-Horizon Forecasting with Elastic Time-Series Transformer

Jiawen Zhang*
DSA, HKUST(GZ)
Guangzhou, China
jiawe.zh@gmail.com

Shun Zheng[†]
Microsoft Research Asia
Beijing, China
shun.zheng@microsoft.com

Xumeng Wen
Microsoft Research Asia
Beijing, China
xumengwen@microsoft.com

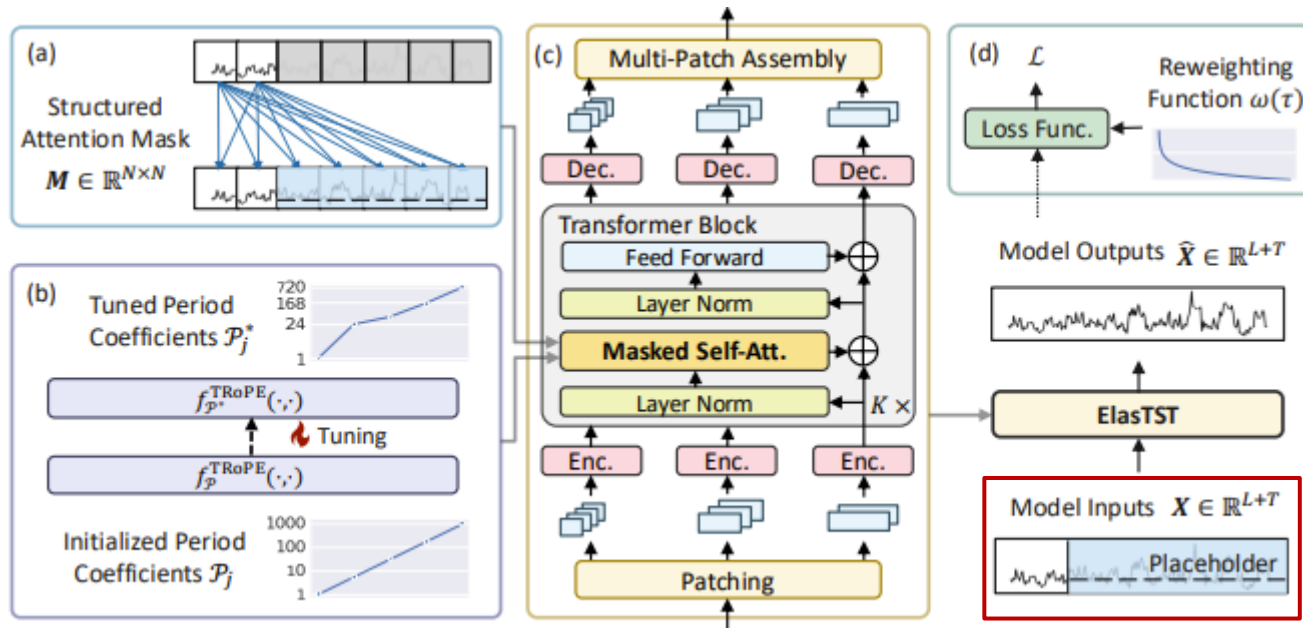
Xiaofang Zhou
CSE, HKUST
Hong Kong SAR, China
zxf@ust.hk

Jiang Bian
Microsoft Research Asia
Beijing, China
jiang.bian@microsoft.com

Jia Li[†]
DSA, HKUST(GZ)
Guangzhou, China
jiale@ust.hk

- arXiv 등록일: 2024-11-04
- 인용 수: 8회(Google Scholar, 2025-12-19)
- Published at NeurIPS 2024
- 기존 시계열 모델들은 특정 예측 길이에 맞춰 학습되는데, 학습 범위를 벗어난 길이를 예측하게 되면 성능이 크게 저하되는 문제가 발생, 이를 해결하기 위해 한 번의 학습으로 다양한 예측 길이에 대해 일관된 성능을 보장하는 ElasTST를 제안

모델 구조



- Model Inputs

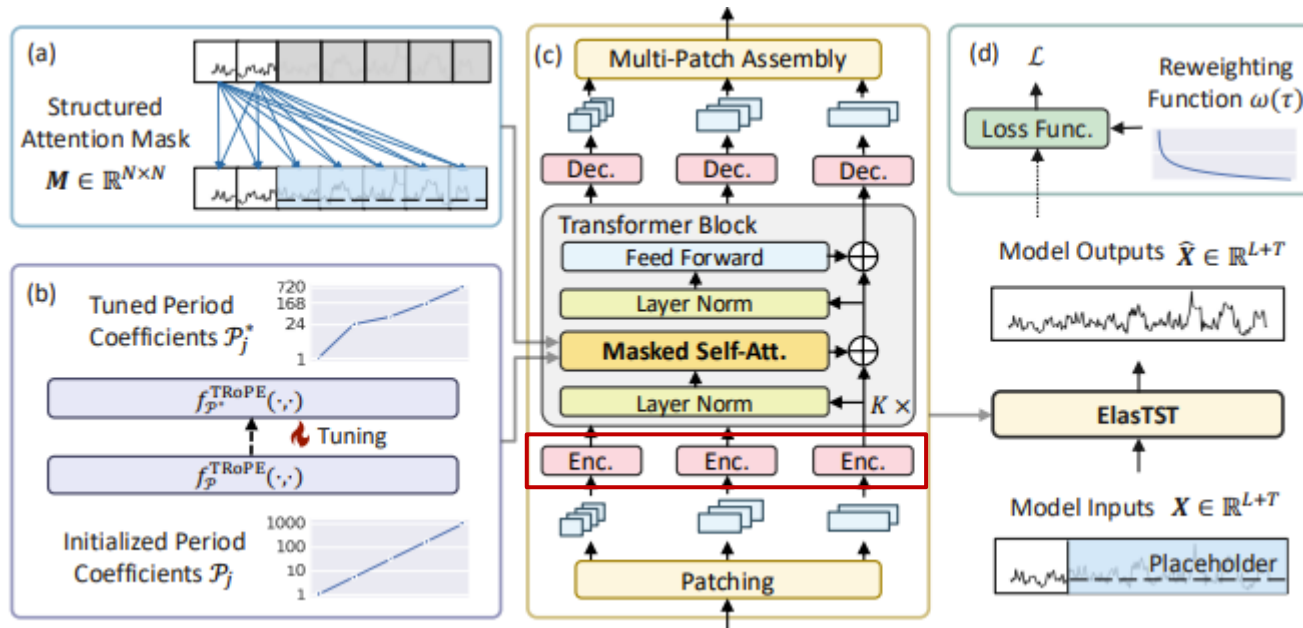
- 과거 시계열 데이터와 예측해야 할 미래 시점인 Placeholder를 Concat하여 입력으로 사용 $X \in \mathbb{R}^{L+T}$

- Multi-Scale Patching

- 입력 X 를 서로 다른 길이의 겹치지 않는 패치 $X^p \in \mathbb{R}^{N+P}$ 로 분할

L : 시퀀스 길이
 T : 예측 길이
 N : 패치 개수
 P : 패치 길이

모델 구조



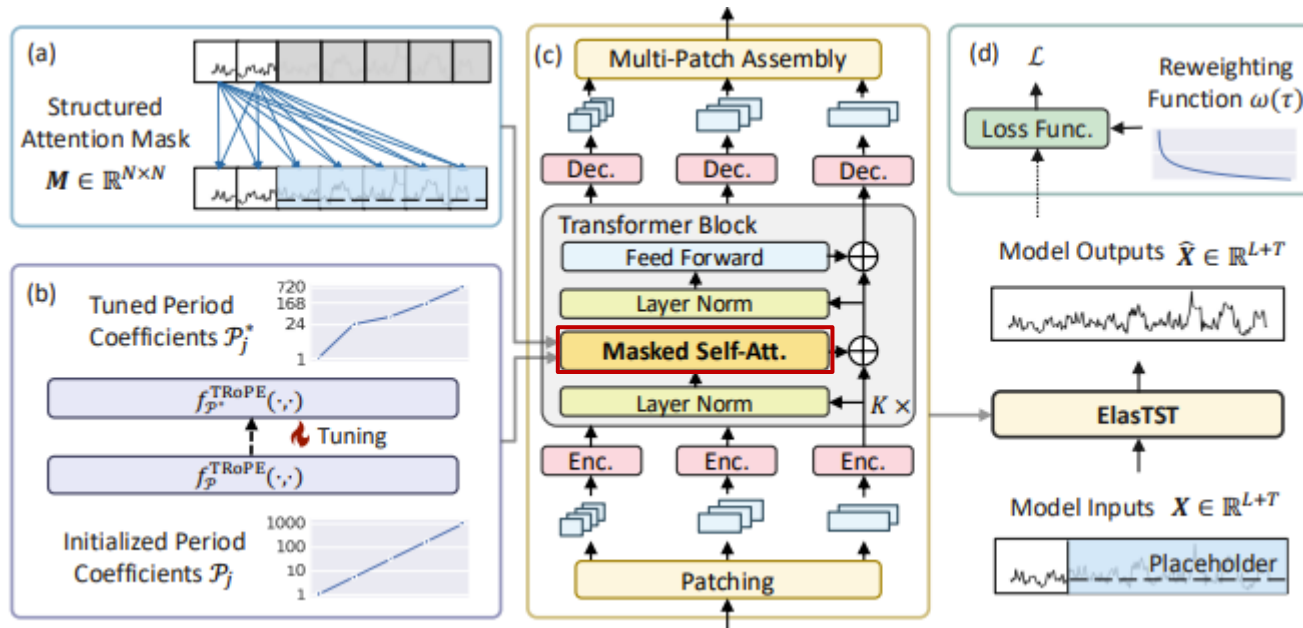
- Encoder

- 분할된 각 패치는 인코더를 통해 D 차원의 임베딩 벡터로 변환 $H = \text{Enc}(X^p)$, $H \in \mathbb{R}^{N+D}$
 - 동일한 D 차원으로 변환해 Transformer Block에서 한 번에 연산

- Layer Norm

- 각 패치별로 D 개의 특징 값들에 대해 동시에 정규화하여 학습 안정화 및 성능 향상

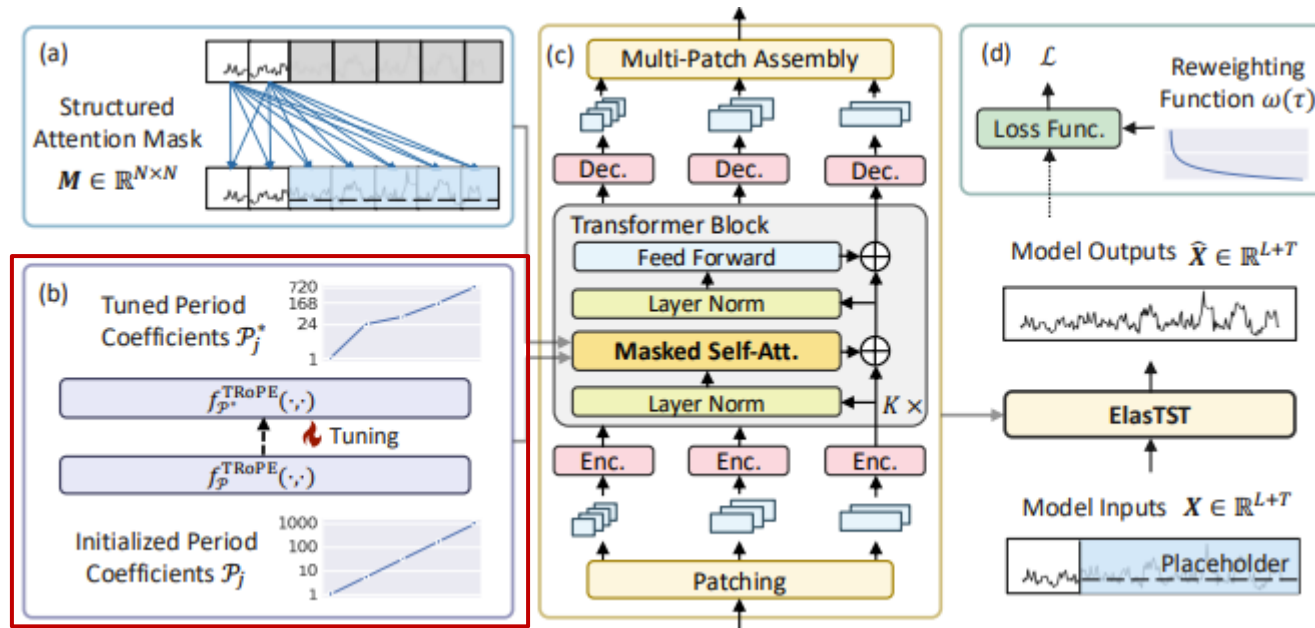
모델 구조



- Masked Self-Attention

- 정규화된 벡터에 가중치 행렬 W^q, W^k, W^v 을 곱해 Q, K, V를 생성하고, TRoPE가 적용된 Q, K의 내적을 통해 Attention Score 계산

모델 구조



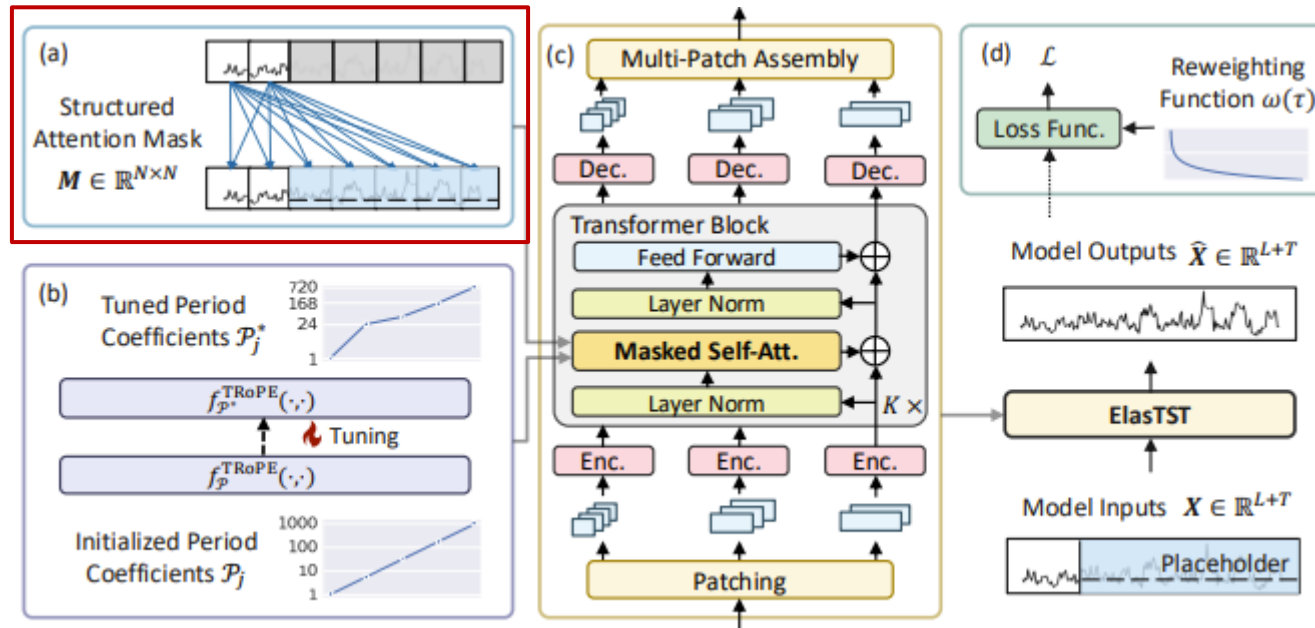
• TRoPE (Tunable Rotary Position Embedding)

- Q 와 K 벡터의 D 개의 숫자들을 각각 가져와 2개씩 짝지어 $D/2$ 개의 쌍으로 묶고, 이후 각각의 쌍은 독립적인 2차원 평면 위의 좌표가 되어 주기를 개별적으로 관리
- 각 평면마다 데이터의 순서 n 과 학습된 주기를 조합한 각도만큼 벡터를 회전 $\theta_j = \frac{2\pi}{P_j} \times n$
- 회전된 Q, K 의 내적을 통해 두 시점 간의 상대적인 각도 차이인 위치 정보를 Attention score에 주입

$$\begin{pmatrix} x'_{2j-1} \\ x'_{2j} \end{pmatrix} = \begin{pmatrix} \cos \theta_j & -\sin \theta_j \\ \sin \theta_j & \cos \theta_j \end{pmatrix} \begin{pmatrix} x_{2j-1} \\ x_{2j} \end{pmatrix}$$

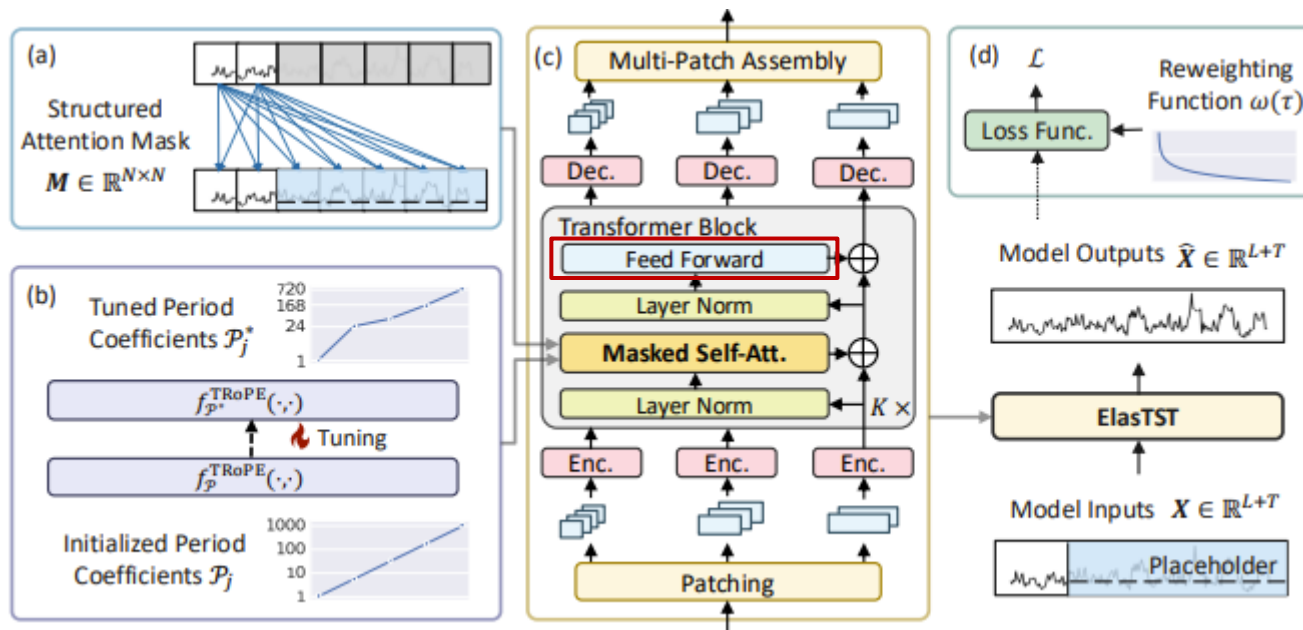
n : 데이터 순서
 θ_j : 최종 각도
 P_j : 주기
 2π : 360도

모델 구조



- Structured Attention Mask
 - Attention Score에 Structured Attention Mask 적용
 - 모든 토큰이 위쪽의 과거 데이터 패치만 참조하여 예측 길이가 바뀌더라도 예측 결과가 변하지 않는 예측 기간 불변성(Horizon-invariant property)을 보장

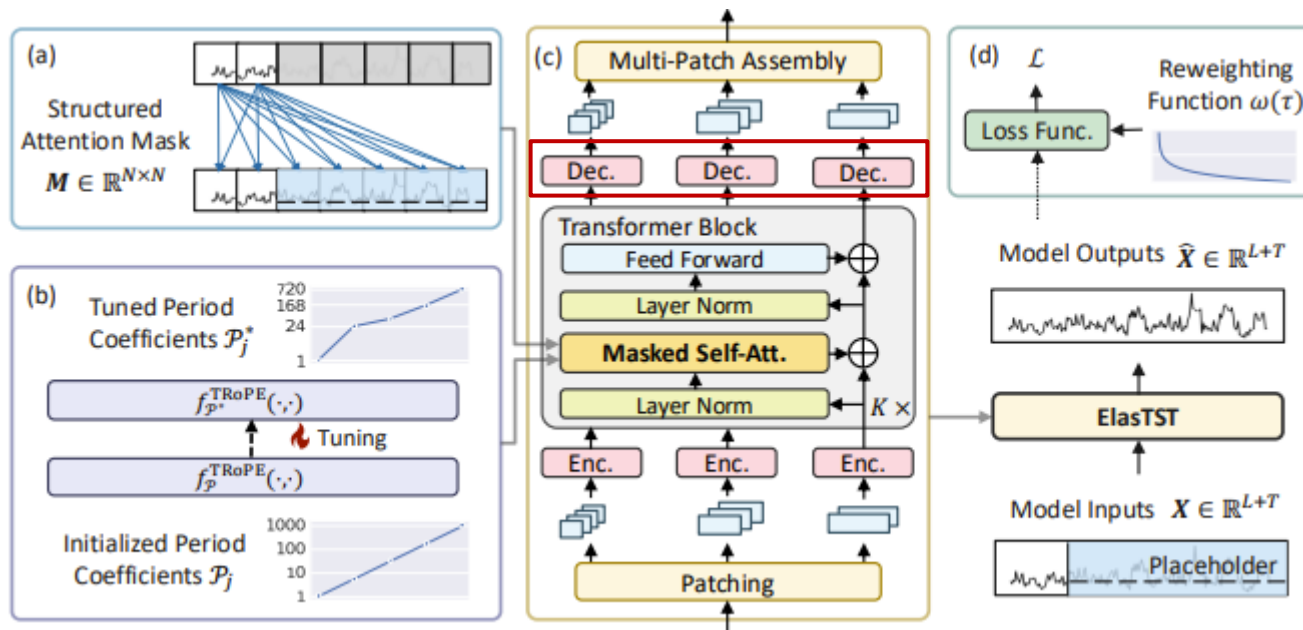
모델 구조



- Feed Forward

- 비선형성 추가, 또한 각 패치마다 동일한 가중치를 공유하며 독립적으로 적용

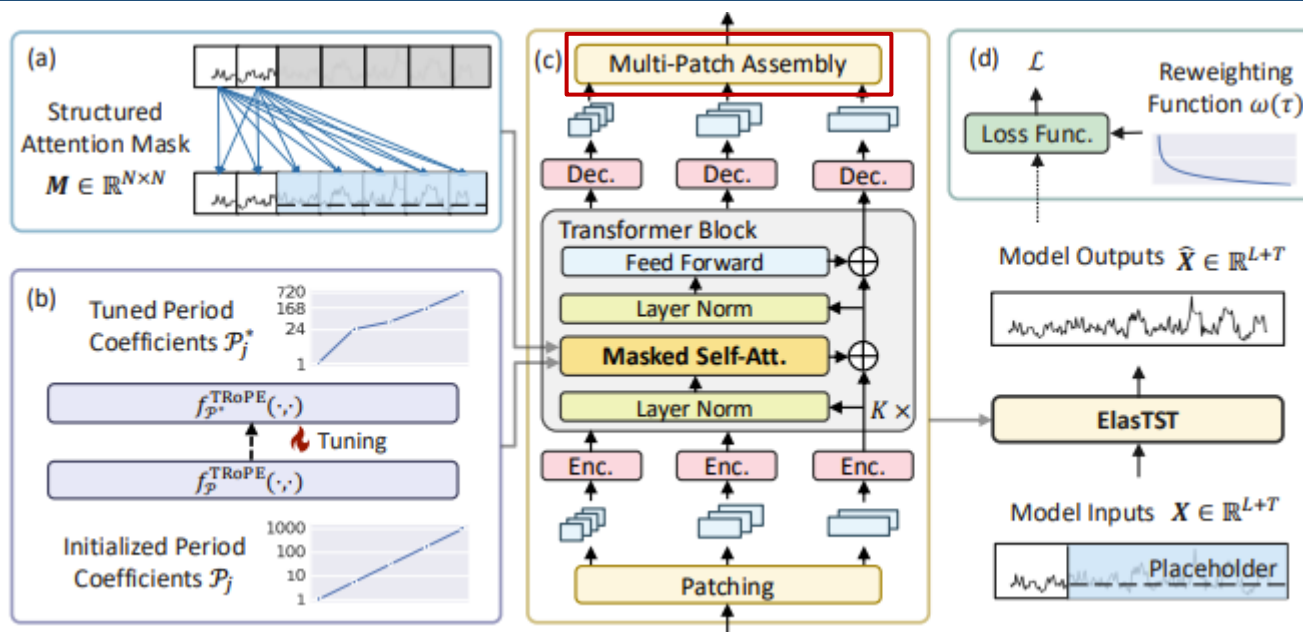
모델 구조



- Decoder

- Transformer Block에서 나온 결과물인 D 차원 특징 벡터를 입력으로 사용하여 원래 패치 길이인 p 차원으로 다시 변환

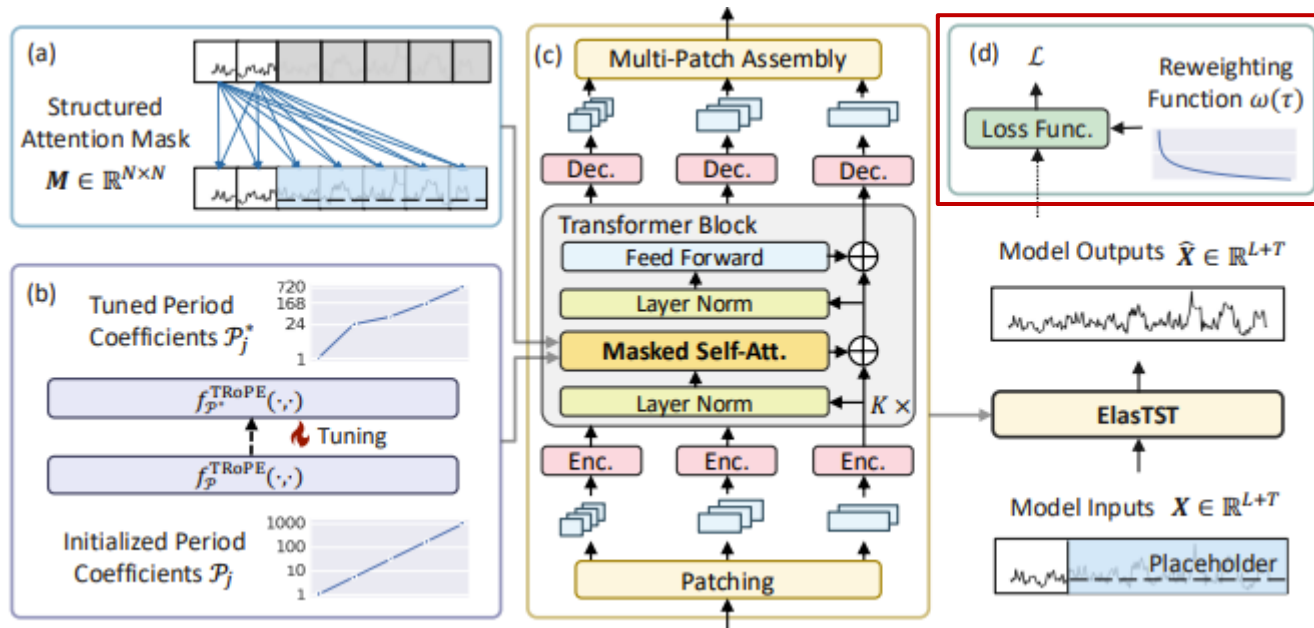
모델 구조



Multi-Patch Assembly

- P 차원으로 복원된 패치들을 이어 붙여 전체 예측 길이만큼의 시퀀스를 만듦
- 이후에 서로 다른 시퀀스를 모두 더한 뒤 시퀀스 개수로 나눠 평균값 출력, 평균값이 최종 예측 결과

모델 구조



• Training Horizon Reweighting

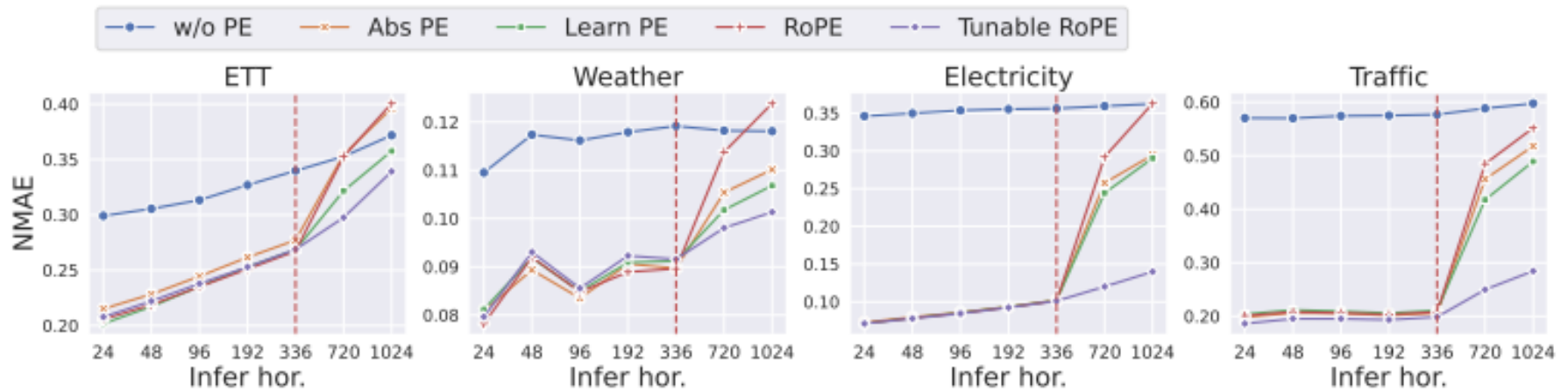
- 매번 예측 길이를 다르게 샘플링하는 대신, 항상 최대치인 T_{max} 까지 학습하도록 하고, 손실 함수를 계산할 때 $[1, T_{max}]$ 범위의 각 시점 τ 에 대해 Reweighting을 적용해 매번 길이를 다르게 샘플링하여 학습한 것과 같은 효과를 냄

주요 모델 성능 비교

	pred len	ElasTST		iTransformer		PatchTST		DLinear		Autoformer	
		NMAE	NRMSE	NMAE	NRMSE	NMAE	NRMSE	NMAE	NRMSE	NMAE	NRMSE
ETTm1	96	0.273. ₀₀₀	0.488.₀₀₀	0.271.₀₀₀	0.568. ₀₀₀	<u>0.272.₀₀₁</u>	<u>0.565.₀₀₁</u>	0.282. ₀₀₂	0.573. ₀₀₁	0.388. ₀₀₁	0.711. ₀₀₃
	192	0.289.₀₀₀	0.520.₀₀₀	0.301. ₀₀₀	0.614. ₀₀₀	<u>0.295.₀₀₁</u>	<u>0.602.₀₀₅</u>	0.309. ₀₀₄	0.617. ₀₀₃	0.442. ₀₀₁	0.820. ₀₀₃
	336	0.314.₀₀₀	0.575.₀₀₀	0.333. ₀₀₀	0.668. ₀₀₀	<u>0.323.₀₀₁</u>	<u>0.645.₀₀₃</u>	0.338. ₀₀₈	0.654. ₀₀₇	0.429. ₀₀₀	0.774. ₀₀₁
	720	0.346.₀₀₀	0.645.₀₀₀	0.376. ₀₀₀	0.741. ₀₀₀	<u>0.353.₀₀₁</u>	<u>0.700.₀₀₅</u>	0.387. ₀₀₆	0.737. ₀₀₅	0.440. ₀₀₀	0.793. ₀₀₀
ETTm2	96	0.150. ₀₀₀	0.227. ₀₀₀	<u>0.137.₀₀₀</u>	0.227. ₀₀₀	0.132.₀₀₁	0.220.₀₀₂	0.138. ₀₀₀	<u>0.226.₀₀₀</u>	0.158. ₀₀₀	0.254. ₀₀₀
	192	0.174. ₀₀₀	<u>0.264.₀₀₀</u>	<u>0.161.₀₀₁</u>	0.266. ₀₀₀	0.157.₀₀₁	0.259.₀₀₂	0.163. ₀₀₃	<u>0.264.₀₀₁</u>	0.175. ₀₀₀	0.283. ₀₀₀
	336	0.191. ₀₀₀	<u>0.289.₀₀₀</u>	<u>0.180.₀₀₀</u>	0.293. ₀₀₀	0.176.₀₀₀	0.286.₀₀₀	0.188. ₀₀₁	0.291. ₀₀₂	0.191. ₀₀₀	0.307. ₀₀₀
	720	<u>0.211.₀₀₀</u>	0.318.₀₀₀	<u>0.211.₀₀₀</u>	0.330. ₀₀₀	0.205.₀₀₁	<u>0.324.₀₀₂</u>	0.219. ₀₀₃	0.327. ₀₀₂	0.217. ₀₀₀	0.338. ₀₀₀
ETTh1	96	0.342. ₀₀₀	0.619.₀₀₀	0.321.₀₀₀	0.626. ₀₀₀	<u>0.328.₀₀₃</u>	0.640. ₀₀₂	0.352. ₀₁₁	0.668. ₀₁₂	0.367. ₀₀₀	0.656. ₀₀₀
	192	<u>0.364.₀₀₀</u>	0.661.₀₀₀	0.359.₀₀₂	<u>0.690.₀₀₀</u>	0.359.₀₀₂	0.705. ₀₀₁	0.393. ₀₀₁	0.745. ₀₀₃	0.392. ₀₀₀	0.706. ₀₀₀
	336	0.371.₀₀₀	0.666.₀₀₀	0.388. ₀₀₀	0.723. ₀₀₀	<u>0.384.₀₀₂</u>	0.740. ₀₀₄	0.419. ₀₀₇	0.778. ₀₀₉	0.398. ₀₀₀	0.711. ₀₀₀
	720	0.376.₀₀₀	0.679.₀₀₀	0.408. ₀₀₀	<u>0.735.₀₀₀</u>	<u>0.397.₀₀₂</u>	0.738. ₀₀₁	0.502. ₀₂₉	0.860. ₀₄₉	0.433. ₀₀₀	0.739. ₀₀₀
ETTh2	96	0.158.₀₀₀	0.239.₀₀₀	<u>0.177.₀₀₀</u>	<u>0.279.₀₀₀</u>	<u>0.177.₀₀₀</u>	0.281. ₀₀₁	0.211. ₀₂₇	0.320. ₀₃₃	0.203. ₀₀₀	0.317. ₀₀₀
	192	0.170.₀₀₀	0.259.₀₀₀	0.203. ₀₀₀	<u>0.314.₀₀₀</u>	<u>0.201.₀₀₁</u>	0.314. ₀₀₁	0.238. ₀₂₈	0.353. ₀₃₀	0.226. ₀₀₀	0.346. ₀₀₀
	336	0.188.₀₀₀	0.282.₀₀₀	0.243. ₀₀₀	0.372. ₀₀₀	<u>0.240.₀₀₁</u>	<u>0.366.₀₀₁</u>	0.284. ₀₀₈	0.407. ₀₁₃	0.264. ₀₀₀	0.398. ₀₀₀
	720	0.215.₀₀₀	0.319.₀₀₀	0.264. ₀₀₀	0.386. ₀₀₀	<u>0.252.₀₀₀</u>	<u>0.371.₀₀₀</u>	0.307. ₀₀₀	0.426. ₀₀₇	0.287. ₀₀₀	0.416. ₀₀₀
Electricity	96	0.085.₀₀₀	<u>0.777.₀₀₀</u>	0.098. ₀₀₀	0.772.₀₀₀	<u>0.086.₀₀₁</u>	0.816. ₀₀₅	0.090. ₀₀₁	0.863. ₀₀₂	0.140. ₀₀₀	0.977. ₀₁₆
	192	<u>0.093.₀₀₀</u>	<u>0.933.₀₀₀</u>	0.106. ₀₀₀	0.916.₀₀₀	0.092.₀₀₁	0.942. ₀₀₇	0.095. ₀₀₁	0.974. ₀₀₁	0.136. ₀₀₀	1.017. ₀₀₀
	336	<u>0.101.₀₀₀</u>	1.063. ₀₀₀	0.115. ₀₀₀	0.985.₀₀₁	0.100.₀₀₀	1.035. ₀₀₃	0.104. ₀₀₀	1.066. ₀₀₄	0.147. ₀₀₀	1.080. ₀₀₆
	720	<u>0.117.₀₀₀</u>	1.289. ₀₀₀	0.133. ₀₀₀	1.110.₀₀₁	0.116.₀₀₀	<u>1.213.₀₀₃</u>	0.122. ₀₀₁	1.259. ₀₀₉	0.159. ₀₀₀	1.283. ₀₀₅
Traffic	96	0.195.₀₀₀	0.461.₀₀₀	<u>0.246.₀₀₀</u>	<u>0.511.₀₀₀</u>	0.248. ₀₀₁	0.527. ₀₀₁	0.356. ₀₀₉	0.645. ₀₁₇	0.293. ₀₀₀	0.560. ₀₀₀
	192	0.193.₀₀₀	0.459.₀₀₀	0.259. ₀₀₀	0.543. ₀₀₀	<u>0.245.₀₀₁</u>	<u>0.528.₀₀₁</u>	0.346. ₀₀₉	0.628. ₀₀₉	0.318. ₀₀₀	0.594. ₀₀₀
	336	0.199.₀₀₀	0.468.₀₀₀	0.283. ₀₀₀	0.571. ₀₀₀	<u>0.257.₀₀₂</u>	<u>0.550.₀₀₁</u>	0.350. ₀₀₈	0.631. ₀₀₈	0.332. ₀₀₀	0.630. ₀₀₀
	720	0.218.₀₀₀	0.497.₀₀₀	0.275. ₀₀₀	0.563. ₀₀₀	<u>0.266.₀₀₁</u>	<u>0.559.₀₀₁</u>	0.365. ₀₀₉	0.659. ₀₀₉	0.341. ₀₀₃	0.611. ₀₀₂
Weather	96	0.086.₀₀₀	0.287.₀₀₀	0.089. ₀₀₀	0.295. ₀₀₀	<u>0.087.₀₀₂</u>	<u>0.294.₀₀₂</u>	0.112. ₀₀₁	0.316. ₀₀₀	0.239. ₀₀₄	0.614. ₀₂₃
	192	<u>0.092.₀₀₀</u>	<u>0.312.₀₀₀</u>	0.093. ₀₀₀	0.299.₀₀₀	<u>0.090.₀₀₁</u>	0.299.₀₀₁	0.122. ₀₀₁	0.331. ₀₀₀	0.213. ₀₀₀	0.533. ₀₀₂
	336	0.091.₀₀₀	<u>0.307.₀₀₀</u>	0.096. ₀₀₀	0.297.₀₀₀	<u>0.092.₀₀₂</u>	0.297.₀₀₁	0.130. ₀₀₂	0.340. ₀₀₂	0.176. ₀₀₀	0.413. ₀₀₁
	720	0.093.₀₀₀	<u>0.308.₀₀₀</u>	0.099. ₀₀₀	0.298.₀₀₀	<u>0.094.₀₀₁</u>	0.298.₀₀₃	0.144. ₀₀₁	0.358. ₀₀₂	0.170. ₀₀₁	0.434. ₀₁₀
Exchange	96	0.026. ₀₀₀	0.039. ₀₀₀	0.025. ₀₀₀	0.039. ₀₀₀	0.023.₀₀₀	0.036.₀₀₀	<u>0.024.₀₀₀</u>	<u>0.037.₀₀₀</u>	0.032. ₀₀₀	0.049. ₀₀₀
	192	0.033.₀₀₀	0.050.₀₀₀	0.036. ₀₀₀	0.056. ₀₀₀	<u>0.034.₀₀₀</u>	<u>0.054.₀₀₀</u>	0.035. ₀₀₀	0.055. ₀₀₀	0.041. ₀₀₀	0.065. ₀₀₀
	336	0.041.₀₀₀	0.062.₀₀₀	0.048. ₀₀₀	0.072. ₀₀₀	<u>0.048.₀₀₀</u>	<u>0.076.₀₀₀</u>	0.048. ₀₀₀	0.072. ₀₀₁	0.056. ₀₀₀	0.091. ₀₀₀
	720	0.059.₀₀₀	0.089.₀₀₀	0.076. ₀₀₀	0.114. ₀₀₀	<u>0.072.₀₀₀</u>	<u>0.106.₀₀₁</u>	0.075. ₀₀₂	0.118. ₀₀₄	0.112. ₀₀₂	0.164. ₀₀₄

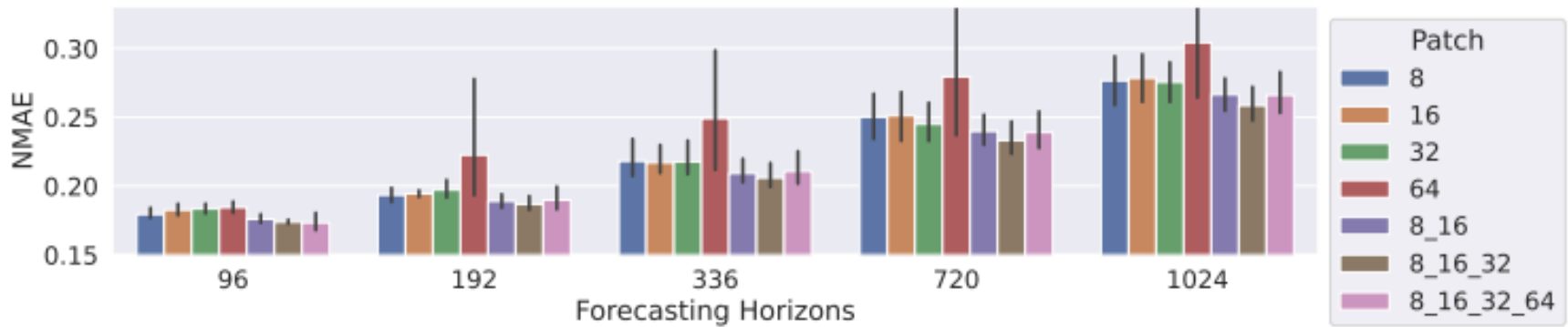
- 기존 Transformer 기반 모델들과 비교하였을 때, ElasTST가 전체적으로 더 뛰어난 성능을 보임

위치 임베딩 비교



- 데이터의 패턴에 맞춰 학습되는 TRoPE로 인해 학습하지 않은 긴 시계열 데이터가 들어와도 안정적으로 예측이 가능

Patch 크기별 성능 비교



- 단일 패치를 사용하는 것보다 여러 패치를 조합하여 사용하는 것이 모든 예측 구간에서 가장 낮은 오차를 기록
- 세밀한 변동과 전역적인 추세를 동시에 포착하는 Multi-Scale이 예측 정확도를 안정화함

실험 세팅

- 사용한 모델: ElasTST
- 재현 실험 데이터셋: ETTh1

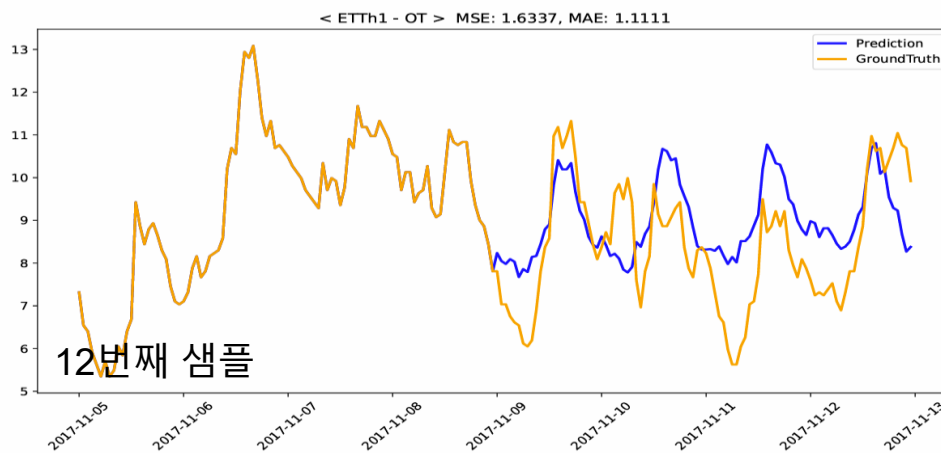
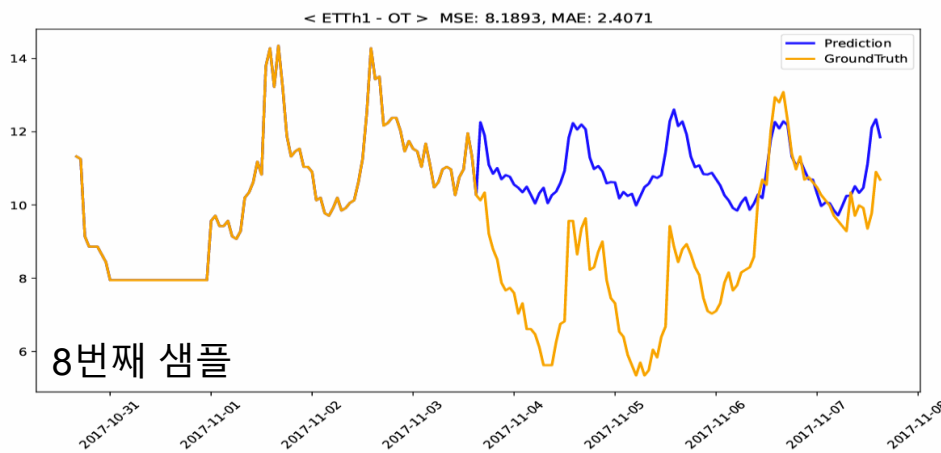
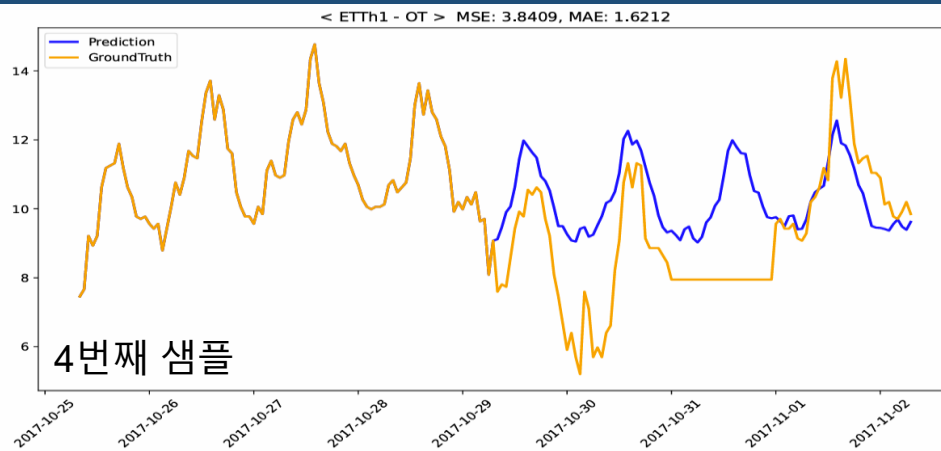
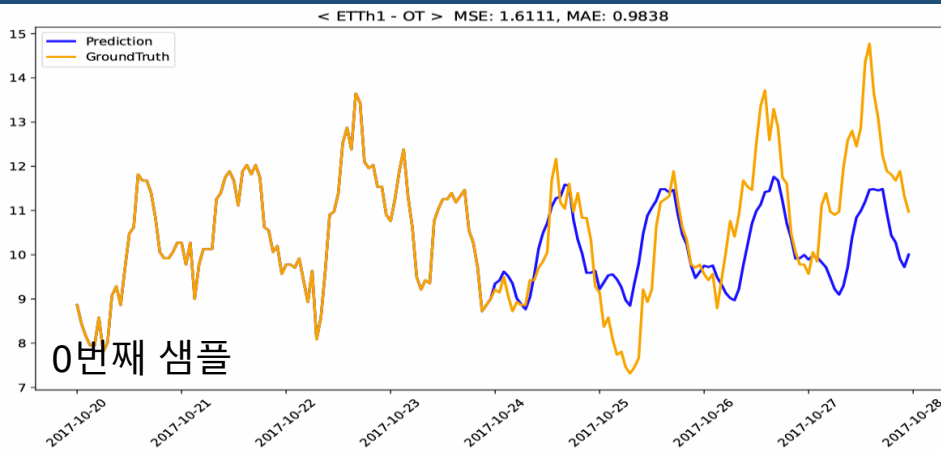
Experiment	ETTh1
Learning rate	10^{-3}
Epoch	20
Batch size	32
Loss function	L1 Loss
Seq_len	96
Pred_len	96/192/336/720
d_model	512
d_ff	512
Patch length	8, 16, 32

ElaSTST 재현 실험 (ETTh1)

- 예측 길이가 길어질수록 어느 정도의 오차가 발생하였지만, 전체적으로 논문과 비슷한 수치가 나온 것을 확인

	ETTh1 Paper		ETTh1 Reproduction	
Pred len	NMAE	NRMSE	NMAE	NRMSE
96	0.342	0.619	0.328	0.627
192	0.364	0.661	0.353	0.672
336	0.371	0.666	0.371	0.704
720	0.376	0.679	0.385	0.698

재현 실험 시각화 (ETTh1)



Seq_len → 96

Pred_len → 96

실험 결과 정리

- 재현 실험
 - 예측 길이가 길어질수록 오차가 발생하였지만 논문과 비슷한 수치가 나옴
- 시각화
 - 전반적으로 예측값이 정답값을 잘 따라가고 있으나, 급격한 변동이나 세밀한 피크 부분은 제대로 반영하지 못함