
What if We Enrich *day-ahead* Solar Irradiance Time Series Forecasting with Spatio-Temporal Context? - *Supplementary material*

Anonymous Author(s)

Affiliation

Address

email

1 Contents

2	1 Experimental details	1
3	1.1 Training setup	1
4	1.1.1 CrossViViT	1
5	1.1.2 Time-series baselines	4
6	2 Additional results and visualisations	5
7	2.1 Performance on (unobserved) validation PAY station, on validation and test years .	5
8	2.2 Visualisations for day ahead time series predictions	6
9	3 Vision Transformers (ViT) and Video Vision Transformers (ViViT)	6

10 1 Experimental details

11 In this section, we present details about training CrossViViT and its Multi-Quantile variant, as well
12 as the time-series baselines.

13 1.1 Training setup

14 1.1.1 CrossViViT

15 CrossViViT integrates two modalities: satellite video data and time-series station data. For both
16 modalities, essential information such as geographic coordinates, elevation, and precise time-stamps
17 is available. In this section, we provide a comprehensive explanation of the encoding process for each
18 feature and conclude by presenting the hyperparameters of the model.

19 We first start by encoding the timestamps. For each time point, we have access to the following *time*
20 *features*: The year, the month, the day, the hour and the minute at which the measurement was made.
21 We use a cyclical embedding to encode these time features discarding the year. For a time feature x ,
22 its corresponding embedding can be expressed as:

$$\left[\sin\left(\frac{2\pi x}{\omega(x)}\right), \cos\left(\frac{2\pi x}{\omega(x)}\right) \right] \quad (1)$$

23 Where $\omega(x)$ is the frequency for time feature x (see Table 1) for the frequency of each time feature).
 24 We concatenate these features to form our final time embedding that we simply concatenate to the
 25 context and the time-series channels respectively. Furthermore, we incorporate the elevation data
 26 for each coordinate in both the context and the time-series. Specifically, the elevation values are
 27 concatenated to their corresponding channels in the context and time-series representations.

28 Regarding the geographic coordinates, we possess information regarding the latitude and longitude
 29 for both the context and the station. These coordinates are normalized so as to lie in $[-1, 1]$:

$$\begin{cases} \text{lat} \leftarrow 2 \left(\frac{\text{lat}+90}{180} \right) - 1 \\ \text{lon} \leftarrow 2 \left(\frac{\text{lon}+180}{360} \right) - 1 \end{cases} \quad (2)$$

30 **ROtary Positional Encoding** Next, we embed these coordinates using ROtary Positional Embed-
 31 ding (ROPE) that we provide a PyTorch implementation for:

```

ROtary Positional Encoding
1 class AxialRotaryEmbedding(nn.Module):
2     def __init__(self, dim, max_freq):
3         super().__init__()
4         self.dim = dim
5         scales = torch.linspace(1.0, max_freq / 2, dim // 4)
6
7         self.register_buffer("scales", scales)
8
9     def forward(self, coords: torch.Tensor):
10        """
11        Args:
12            coords (torch.Tensor): Coordinates of shape [B, 2, height, width]
13        """
14        seq_x = coords[:, 0, 0, :]
15        seq_x = seq_x.unsqueeze(-1)
16        seq_y = coords[:, 1, :, 0]
17        seq_y = seq_y.unsqueeze(-1)
18
19        scales = self.scales[ *((None, None)), Ellipsis]
20        scales = scales.to(coords)
21
22        seq_x = seq_x * scales * pi
23        seq_y = seq_y * scales * pi
24
25        x_sinu = repeat(seq_x, "b i d -> b i j d", j=seq_y.shape[1])
26        y_sinu = repeat(seq_y, "b j d -> b i j d", i=seq_x.shape[1])
27
28        sin = torch.cat((x_sinu.sin(), y_sinu.sin()), dim=-1)
29        cos = torch.cat((x_sinu.cos(), y_sinu.cos()), dim=-1)
30
31        sin, cos = map(lambda t: rearrange(t, "b i j d -> b (i j) d"), (sin,
32        ↪ cos))
33        sin, cos = map(lambda t: repeat(t, "b n d -> b n (d j)", j=2), (sin,
34        ↪ cos))
35        return sin, cos
  
```

Table 1: Frequency of each time feature.

Time feature	Frequency
Month	12
Day	31
Hour	24
Minute	60

32 **Training configuration** CrossViViT was trained on two RTX8000 GPUs, over 17 epochs with
33 early stopping. Its Multi-Quantile variant was also trained on two RTX8000 GPUs, over 12 epochs
34 with early stopping. The remaining settings are identical for both variants: An effective batch size of
35 20 was utilized for both models; The training process employed the AdamW optimizer (Loshchilov
36 and Hutter, 2019) with a weight decay of 0.05; A cosine warmup strategy was implemented, gradually
37 increasing the learning rate from 0 to 0.0016 over five epochs before starting the decay phase.

38 Below, we highlight the relevant model hyperparameters for CrossViViT and Multi-Quantile Cross-
39 ViViT:

CrossViViT hyperparameters

```
1 patch_size: [8, 8]
2 use_glu: True
3 max_freq: 128
4 num_mlp_heads: 1
5
6 ctx_masking_ratio: 0.99
7 ts_masking_ratio: 0
8
9 # These hyperparameters apply to the encoding transformers and cross-attention
10 dim: 384
11 depth: 16
12 heads: 12
13 mlp_ratio: 4
14 dim_head: 64
15 dropout: 0.4
16
17 # These only apply to the decoding transformer
18 decoder_dim: 128
19 decoder_depth: 4
20 decoder_heads: 6
21 decoder_dim_head: 128
```

Multi-Quantile CrossViViT hyperparameters

```
1 patch_size: [8, 8]
2 use_glu: True
3 max_freq: 128
4 num_mlp_heads: 11
5 quantiles: [0.02, 0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9, 0.98]
6
7 ctx_masking_ratio: 0.99
8 ts_masking_ratio: 0
9
10 # These hyperparameters apply to the encoding transformers and cross-attention
11 dim: 256
12 depth: 16
13 heads: 12
14 mlp_ratio: 4
15 dim_head: 64
16 dropout: 0.4
17
18 # These only apply to the decoding transformer
19 decoder_dim: 128
20 decoder_depth: 4
21 decoder_heads: 6
22 decoder_dim_head: 128
```

40 1.1.2 Time-series baselines

41 We conducted training on a total of nine baseline models, and we emphasize the importance of the
42 hyperparameters used for each of these models. Below, we provide a comprehensive overview of the
43 hyperparameters employed in our study:

- 44 • **seq_len**: Input sequence length.
- 45 • **label_len**: Start token length.
- 46 • **pred_len**: Prediction sequence length.
- 47 • **enc_in**: Encoder input size.
- 48 • **dec_in**: Decoder input size.
- 49 • **e_layers**: Number of encoder layers.
- 50 • **d_layers**: Number of decoder layers.
- 51 • **c_out**: Output size.
- 52 • **d_model**: Dimension of model.
- 53 • **n_heads**: Number of attention heads.
- 54 • **d_ff**: Dimension of Fully Connected Network.
- 55 • **factor**: Attention factor.
- 56 • **embed**: Time features encoding, options:[timeF, fixed, learned].
- 57 • **distil**: Whether to use distilling in encoder.
- 58 • **moving average**: Window size of moving average kernel.

59 We adapted the majority of the baselines using the Time Series Library (TSlb (Wu et al., 2023)),
60 which served as a valuable resource in our experimentation. We refer the reader to the original papers,
61 which served as a base for the hyperparameters utilized in our study, in order to have a comprehensive
62 understanding of the models and the training settings.

LightTS	FiLM	DLinear
1 model: 2 enc_in: 10 3 seq_len: 48 4 pred_len: 48 5 d_model: 256 6 dropout: 0.05 7 chunk_size: 24	1 model: 2 enc_in: 10 3 seq_len: 48 4 label_len: 24 5 pred_len: 48 6 e_layers: 2 7 ratio: 0.4	1 model: 2 enc_in: 10 3 seq_len: 48 4 pred_len: 48 5 moving_avg: 25 6 individual: False
Crossformer	Reformer	PatchTST
1 model: 2 enc_in: 10 3 seq_len: 48 4 pred_len: 48 5 d_model: 1024 6 n_heads: 2 7 e_layers: 4 8 d_ff: 2048 9 factor: 10 10 dropout: 0.01	1 model: 2 enc_in: 10 3 c_out: 1 4 seq_len: 48 5 pred_len: 48 6 d_model: 512 7 n_heads: 8 8 e_layers: 3 9 d_ff: 2048 10 factor: 5 11 dropout: 0.05 12 embed: timeF 13 activation: gelu	1 model: 2 enc_in: 10 3 c_out: 1 4 seq_len: 48 5 pred_len: 48 6 d_model: 1024 7 n_heads: 6 8 e_layers: 3 9 d_ff: 2048 10 factor: 10 11 dropout: 0.05 12 embed: timeF 13 activation: gelu

Autoformer		Informer		FEDformer	
1	model:	1	model:	1	model:
2	enc_in: 10	2	enc_in: 10	2	enc_in: 10
3	dec_in: 10	3	dec_in: 10	3	dec_in: 10
4	c_out: 1	4	c_out: 1	4	c_out: 1
5	seq_len: 48	5	label_len: 24	5	seq_len: 48
6	label_len: 24	6	pred_len: 48	6	label_len: 24
7	pred_len: 48	7	d_model: 2048	7	pred_len: 48
8	moving_avg: 25	8	n_heads: 4	8	moving_avg: 25
9	d_model: 1024	9	e_layers: 2	9	d_model: 512
10	n_heads: 8	10	d_layers: 2	10	n_heads: 8
11	e_layers: 3	11	d_ff: 2048	11	e_layers: 3
12	d_layers: 2	12	factor: 5	12	d_layers: 2
13	d_ff: 2048	13	dropout: 0.1	13	d_ff: 2048
14	factor: 10	14	embed: timeF	14	dropout: 0.05
15	dropout: 0.01	15	activation: gelu	15	version: fourier
16	embed: timeF	16	distil: True	16	mode_select: random
17	activation: gelu			17	modes: 32

63 The training of the baselines took place on a single RTX8000 GPU over the course of 100 epochs.
64 During training, a batch size of 64 was consistently employed. For model optimization, we utilized
65 the AdamW optimizer (Loshchilov and Hutter, 2019), incorporating a weight decay value set to 0.05.
66 Moreover, we implemented a learning rate reduction strategy known as Reduce Learning Rate on
67 Plateau, which gradually decreased the learning rate by a factor of 0.5 after a patience of 10 epochs.

68 For the hyperparameter tuning of the baselines, we employed the Orion package (Bouthillier et al.,
69 2022), which is an asynchronous framework designed for black-box function optimization. Its
70 primary objective is to serve as a meta-optimizer specifically tailored for machine learning models.
71 As an example, we present the details of the Crossformer hyperparameter tuning scheme, showcasing
the approach we followed to optimize its performance:

Crossformer hyperparameters tuning	
1	params:
2	optimizer.lr: loguniform(1e-8, 0.1)
3	optimizer.weight_decay: loguniform(1e-10, 1)
4	d_model: choices([128, 256, 512, 1024, 2048])
5	d_ff: choices([1024, 2048, 4096])
6	n_heads: choices([1, 2, 4, 8, 16])
7	e_layers: choices([1, 2, 3, 4, 5])
8	dropout: choices([0.01, 0.05, 0.1, 0.2, 0.25])
9	factor: choices([2, 5, 10])
10	max_epochs: fidelity(low=5, high=100, base=3)

72

73 2 Additional results and visualisations

74 2.1 Performance on (unobserved) validation PAY station, on validation and test years

75 This section presents a comprehensive comparative analysis that assesses the performance of Cross-
76 ViViT in relation to various timeseries approaches and solar irradiance baselines specifically for the
77 PAY station. The evaluation encompasses both the validation period (2017-2019) and the test period
78 (2020-2022), with the validation period serving as the basis for model selection. The results, as
79 presented in Table 2, offer compelling evidence of the superior performance of CrossViViT compared
80 to the alternative approaches across all evaluation splits. These findings underscore the crucial role of
81 accurately capturing cloud dynamics in solar irradiance forecasting, which is particularly pronounced
82 in the "Hard" splits.

Table 2: Comparison of model performances in the **val station** PAY, during **test years** (2020-2022) and **val years** (2017-2019). We report the MAE and RMSE for the easy and difficult splits (presented in the main paper) along with the number of data points for each split. We add the MAE resulting from the Multi-Quantile CrossViViT median prediction, along with p_t , the probability for the ground-truth to be included within the interval, averaged across time steps.

Models	Parameters	PAY (Test years - 2020-2022)						PAY (Val years - 2017-2019)					
		All (12171)		Easy (8053)		Hard (4118)		All (27166)		Easy (16683)		Hard (10483)	
		MAE	RMSE	MAE	RMSE	MAE	RMSE	MAE	RMSE	MAE	RMSE	MAE	RMSE
Persistence	N/A	61.9	140.05	44.98	105.47	94.99	190.31	63.61	141.16	48.57	108.91	87.54	180.99
Fourier ₃	N/A	69.21	129.29	52.17	91.61	102.53	181.63	71.71	130.58	56.03	94.81	96.68	172.86
Fourier ₄	N/A	65.67	131.13	48.77	93.68	98.72	183.48	67.35	132.41	51.8	96.71	92.08	174.79
Fourier ₅	N/A	63.39	132.38	46.14	95.08	97.14	184.71	65.71	133.87	50.03	98.61	90.66	175.98
Clear Sky (Ineichen, 2016)	N/A	65.99	138.78	56.41	116.72	84.72	174.03	62.12	131.69	52.59	109.91	77.28	160.36
ReFormer (Kitaev et al., 2020)	8.6M	59.15	109.75	51.67	91.67	73.79	138.44	59.3	109.3	54.06	94.31	67.64	129.62
Informr (Zhou et al., 2021)	56.7M	79.18	133.47	75.8	126.07	85.78	146.86	77.51	131.72	75.65	126.58	80.46	139.52
FiLM (Zhou et al., 2022b)	9.4M	69.87	124.86	55.93	91.15	97.11	172.72	71.1	125.32	58.35	94.17	91.39	163.06
PatchTST (Nie et al., 2023)	9.6M	63.21	131.64	53.06	113.78	83.06	160.93	63.88	131.14	55.47	115.59	77.25	152.64
LighTS (Zhang et al., 2022)	32K	58.57	111.22	48.93	88.69	77.4	145.53	58.22	110.84	50.72	91.24	70.16	136.33
CrossFormer (Zhang and Yan, 2023)	227M	59.64	111.03	51.04	90.34	76.45	143.08	59.39	111.49	52.55	93.15	70.29	135.66
FEDFormer (Zhou et al., 2022a)	23.6M	63.44	110.62	58.15	97.33	73.79	132.83	62.46	109.92	59.88	100.61	66.57	123.3
DLinear (Zeng et al., 2022)	4.7K	75.09	128.64	59.42	93.82	105.74	178.04	76.78	129.45	62.64	97.93	99.29	167.81
AutoFormer (Wu et al., 2021)	50.4M	73.36	117.22	68.89	105.52	82.12	137.25	71.39	114.96	69.42	106.79	74.54	126.88
CrossViViT	145M	51.47	107.73	41.59	85.14	70.81	141.87	52.02	108.77	44.71	90.47	63.65	132.79
		MAE	p_t	MAE	p_t	MAE	p_t	MAE	p_t	MAE	p_t	MAE	p_t
Multi-Quantile CrossViViT	78.8M	59.06	0.83	54.96	0.83	67.06	0.84	57.65	0.86	51.75	0.87	67.05	0.86

83 2.2 Visualisations for day ahead time series predictions

84 This section presents visualizations of predictions generated by CrossViViT and CrossFormer on the
 85 PAY station for both the validation period (2017-2019) (see Figure 2) and the test period (2020-2022)
 86 (see Figure 1). We also present visualizations of predictions generated by the Multi-Quantile version,
 87 for the two periods, on the PAY station. The predictions depicted in red are the median (50% quantile)
 88 estimation from the model, and the generated prediction interval is defined as the interval between
 89 the two predefined quantiles: $[q_{0.02}, q_{0.98}]$

90 3 Vision Transformers (ViT) and Video Vision Transformers (ViViT)

91 **ViT** The Vision Transformer (ViT) model (Dosovitskiy et al., 2020) leverages self-attention mecha-
 92 nisms inspired by the popular transformer architecture (Bahdanau et al., 2015; Vaswani et al., 2017)
 93 to efficiently process images. The input image of dimensions $H \times W$ is divided into non-overlapping
 94 patches $x_i \in \mathbb{R}^{h \times w}$, which are linearly projected and transformed into d -dimensional vector tokens
 95 $z_i \in \mathbb{R}^d$ using a learned weight matrix \mathbf{E} that applies 2D convolution. The sequence of patches is
 96 defined as $\mathbf{z} = [z_{\text{class}}, \mathbf{E}x_1, \dots, \mathbf{E}x_P] + \mathbf{E}_{\text{pos}}$, where z_{class} is an optional learned classification token
 97 representing the class label, and $\mathbf{E}_{\text{pos}} \in \mathbb{R}^{P \times d}$ is a 1D learned positional embedding that encodes
 98 position information.

99 To extract global features from the image, the embedded patches undergo K transformer layers.
 100 Each transformer layer k consists of Multi-Headed Self-Attention (MSA) (Vaswani et al., 2017),
 101 layer normalization (LN) (Ba et al., 2016), and MLP blocks with residual connections. The MLPs
 102 consist of two linear layers separated by the Gaussian Error Linear Unit (GELU) activation function
 103 (Hendrycks and Gimpel, 2016). The output of the final layer can be used for image classification,
 104 either by directly utilizing the classification token or by applying global average pooling to all tokens
 105 \mathbf{z}^L if the classification token was not used initially.

106 **ViViT** The Video Vision Transformer (Arnab et al., 2021) extends the ViT model to handle video
 107 classification tasks by incorporating both spatial and temporal dimensions within a transformer-like
 108 architecture. The authors propose multiple versions of the model, with a key consideration being how
 109 to embed video clips for attention computation.

110 Two approaches are presented: (1) *Uniform frame sampling* and (2) *Tubelet embedding*. The first
 111 method involves uniformly sampling n_t frames from the video clip, applying the same ViT embedding
 112 method Dosovitskiy et al. (2020) to each 2D frame independently, and concatenating the resulting
 113 tokens into a single sequence. The second method extracts non-overlapping spatio-temporal "tubes"

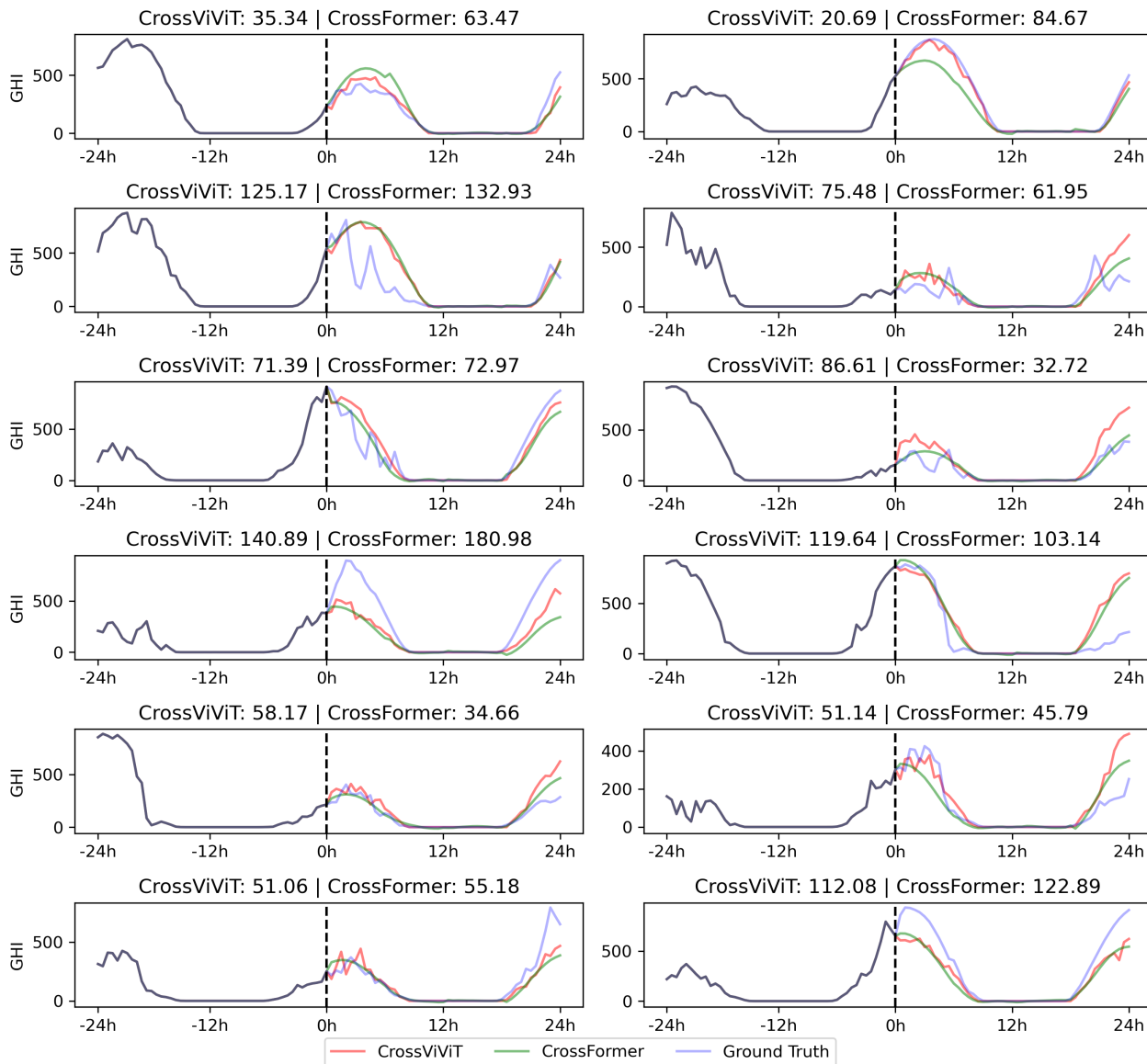


Figure 1: Visualizations of CrossViViT and CrossFormer are presented for a subset of 12 randomly selected days from the "Hard" split on the PAY station during the test period spanning from 2020 to 2022.

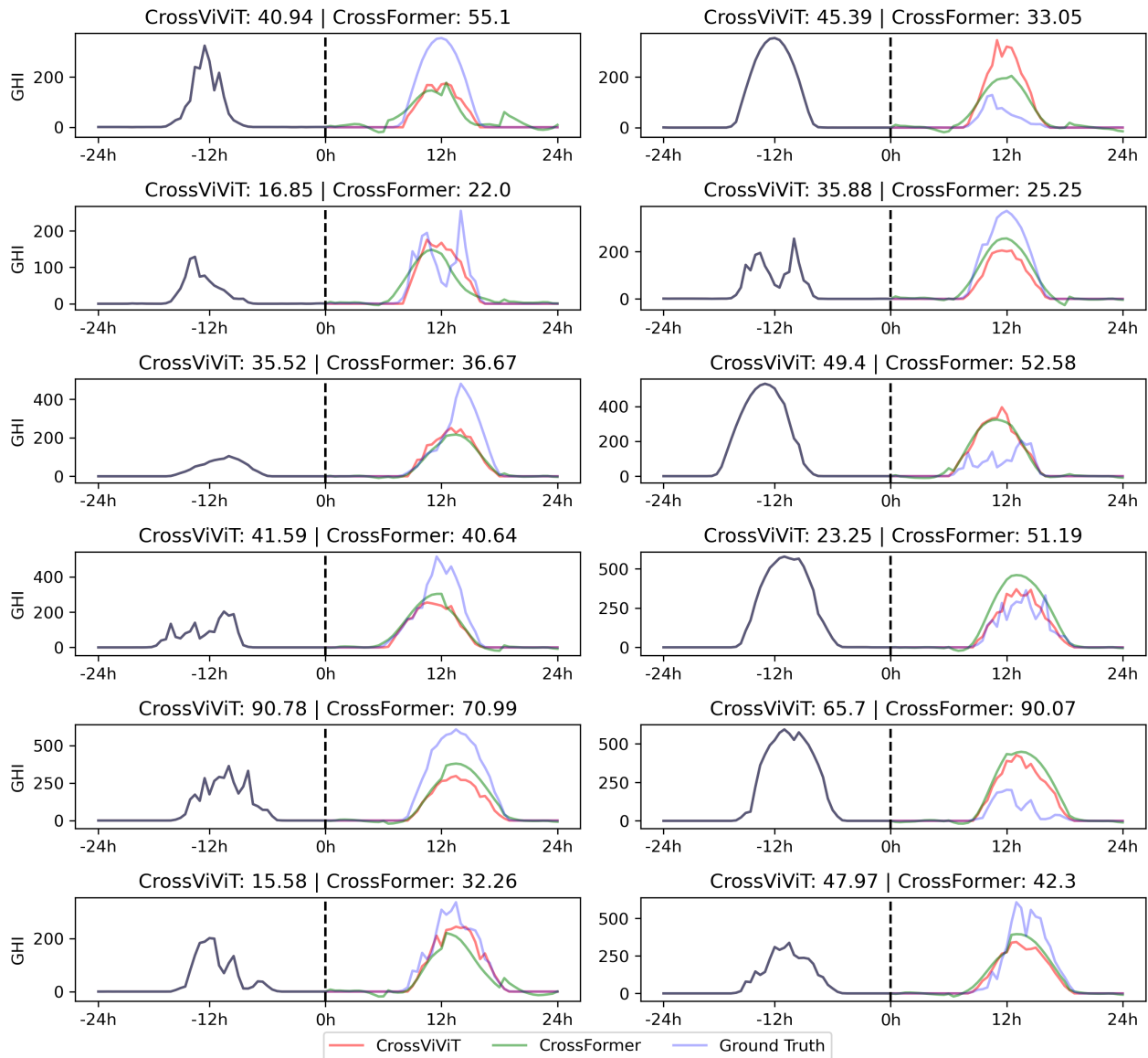


Figure 2: Visualizations of CrossViViT and CrossFormer are presented for a subset of 12 randomly selected days from the "Hard" split on the PAY station during the validation period spanning from 2017 to 2019.

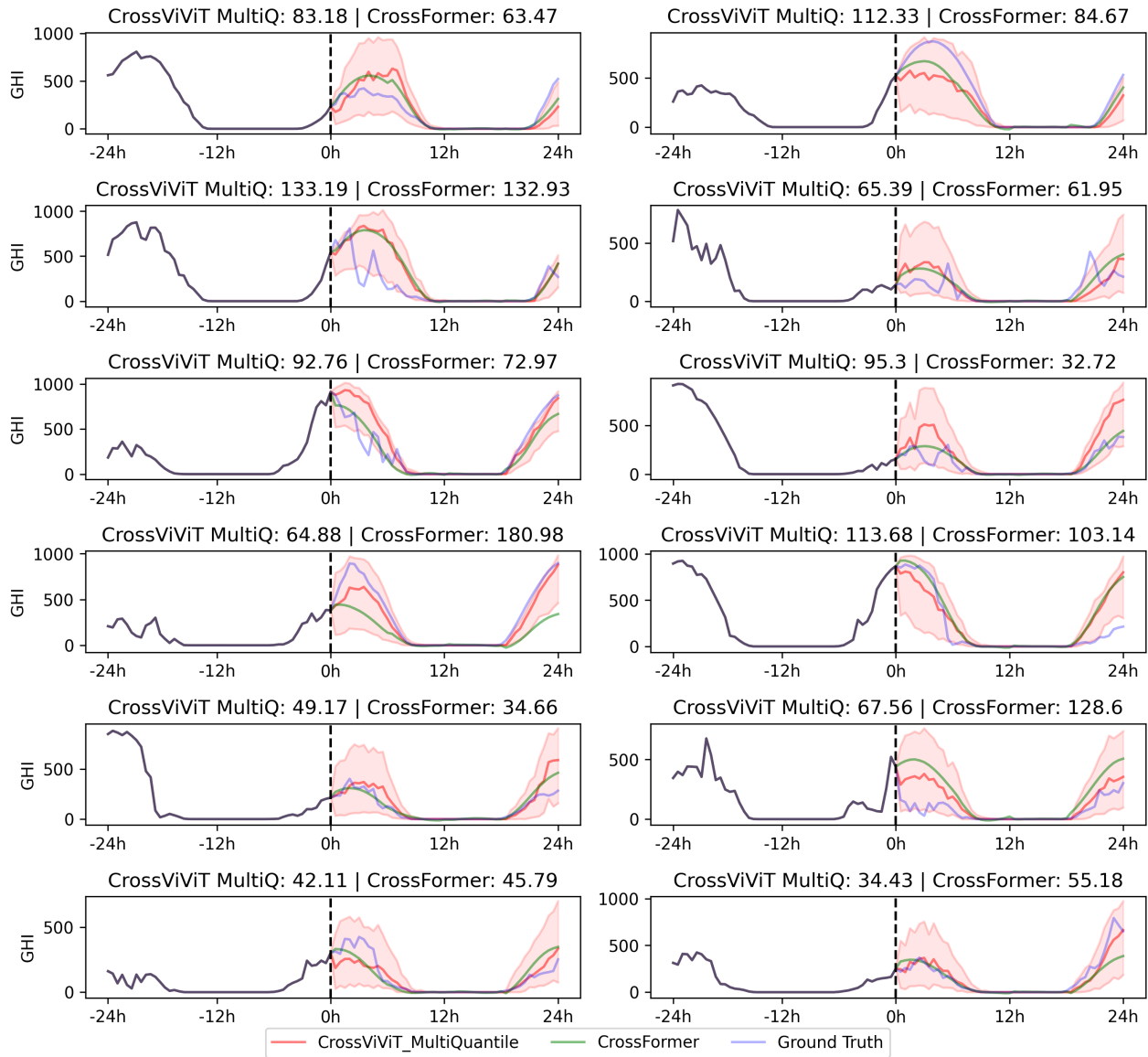


Figure 3: Visualizations of Multi-Quantile CrossViViT (median prediction and $[q_{0.02}, q_{0.98}]$ prediction interval) and CrossFormer predictions are presented for a subset of 12 randomly selected days from the "Hard" split on the PAY station during the validation period spanning from 2020 to 2022.

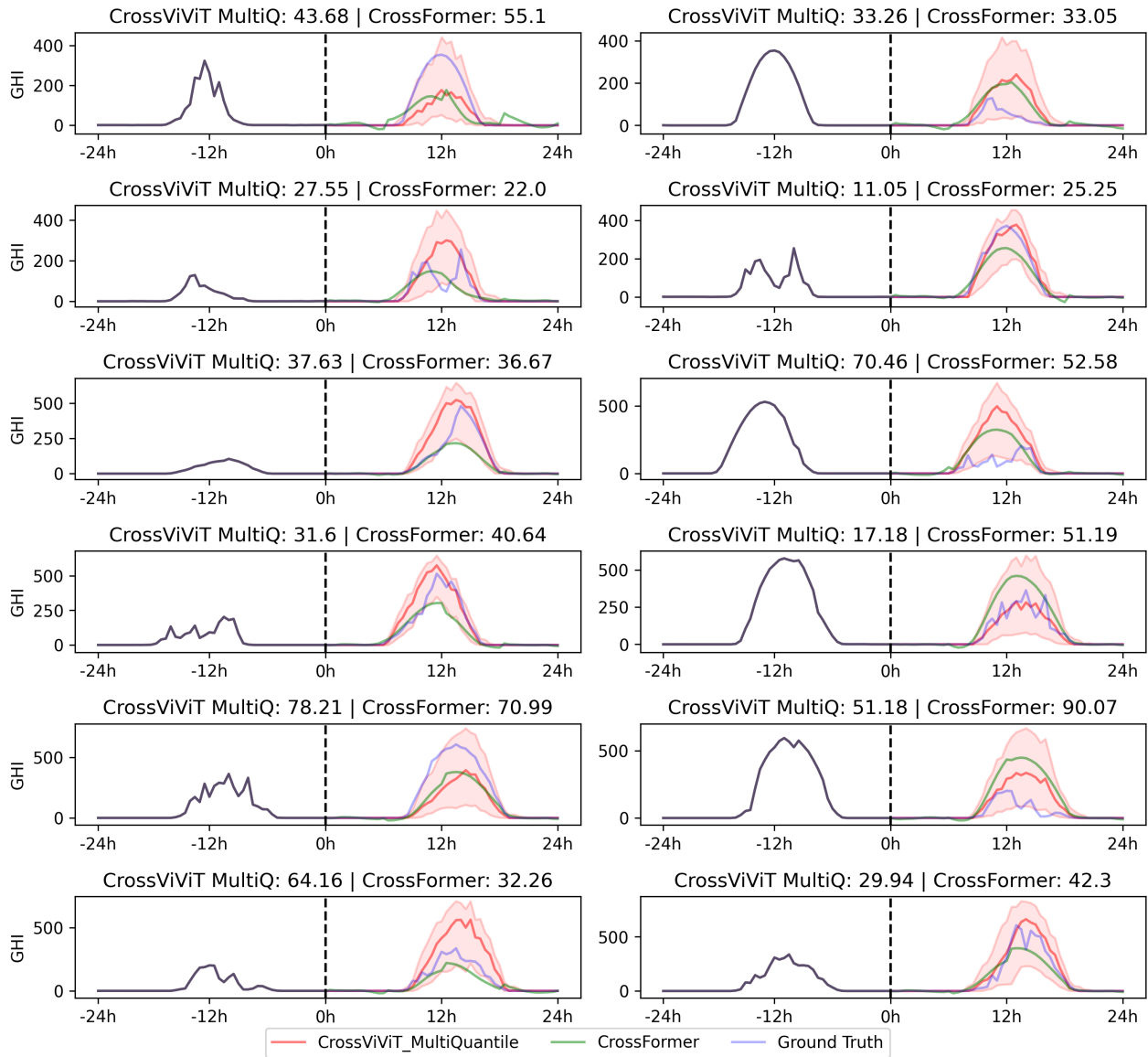


Figure 4: Visualizations of Multi-Quantile CrossViViT (median prediction and $[q_{0.02}, q_{0.98}]$ prediction interval) and CrossFormer predictions are presented for a subset of 12 randomly selected days from the "Hard" split on the PAY station during the validation period spanning from 2017 to 2019.

114 from the input volume and linearly projects them into \mathbb{R}^d , effectively extending ViT’s embedding
115 technique to 3D data, akin to performing a 3D convolution. Intuitively, the second method allows
116 for the fusion of spatio-temporal information during tokenization, whereas the first method requires
117 independent temporal fusion post-tokenization. Experimental results, however, show only slightly
118 superior performance for the second method in specific scenarios, despite its significantly higher
119 computational complexity (Arnab et al., 2021).

120 References

- 121 Arnab, A., Deghani, M., Heigold, G., Sun, C., Lučić, M., and Schmid, C. (2021). Vivit: A video
122 vision transformer. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*
123 (*ICCV*), pages 6836–6846.
- 124 Ba, J. L., Kiros, J. R., and Hinton, G. E. (2016). Layer normalization. *arXiv preprint*
125 *arXiv:1607.06450*.
- 126 Bahdanau, D., Cho, K., and Bengio, Y. (2015). Neural machine translation by jointly learning to align
127 and translate. In Bengio, Y. and LeCun, Y., editors, *3rd International Conference on Learning*
128 *Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*.
- 129 Bouthillier, X., Tsigotakis, C., Corneau-Tremblay, F., Schweizer, T., Dong, L., Delaunay, P., Nor-
130 mandin, F., Bronzi, M., Suhubdy, D., Askari, R., Noukhovitch, M., Xue, C., Ortiz-Gagné, S.,
131 Breuleux, O., Bergeron, A., Bilaniuk, O., Bocco, S., Bertrand, H., Alain, G., Serdyuk, D., Hen-
132 derson, P., Lamblin, P., and Beckham, C. (2022). Epistimio/orion: Asynchronous Distributed
133 Hyperparameter Optimization.
- 134 Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., Deghani, M.,
135 Minderer, M., Heigold, G., Gelly, S., et al. (2020). An image is worth 16x16 words: Transformers
136 for image recognition at scale. *arXiv preprint arXiv:2010.11929*.
- 137 Hendrycks, D. and Gimpel, K. (2016). Gaussian error linear units (gelus). *arXiv preprint*
138 *arXiv:1606.08415*.
- 139 Ineichen, P. (2016). Validation of models that estimate the clear sky global and beam solar irradiance.
140 *Solar Energy*, 132:332–344.
- 141 Kitaev, N., Kaiser, L., and Levskaya, A. (2020). Reformer: The efficient transformer. In *International*
142 *Conference on Learning Representations*.
- 143 Loshchilov, I. and Hutter, F. (2019). Decoupled weight decay regularization. In *7th International*
144 *Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*.
145 OpenReview.net.
- 146 Nie, Y., Nguyen, N. H., Sinthong, P., and Kalagnanam, J. (2023). A time series is worth 64 words:
147 Long-term forecasting with transformers.
- 148 Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, Ł., and
149 Polosukhin, I. (2017). Attention is all you need. *Advances in neural information processing*
150 *systems*, 30.
- 151 Wu, H., Hu, T., Liu, Y., Zhou, H., Wang, J., and Long, M. (2023). Timesnet: Temporal 2d-
152 variation modeling for general time series analysis. In *International Conference on Learning*
153 *Representations*.
- 154 Wu, H., Xu, J., Wang, J., and Long, M. (2021). Autoformer: Decomposition transformers with
155 auto-correlation for long-term series forecasting. *Advances in Neural Information Processing*
156 *Systems*, 34:22419–22430.
- 157 Zeng, A., Chen, M.-H., Zhang, L., and Xu, Q. (2022). Are transformers effective for time series
158 forecasting? *ARXIV.ORG*.
- 159 Zhang, T., Zhang, Y., Cao, W., Bian, J., Yi, X., Zheng, S., and Li, J. (2022). Less is more: Fast
160 multivariate time series forecasting with light sampling-oriented mlp structures.

- 161 Zhang, Y. and Yan, J. (2023). Crossformer: Transformer utilizing cross-dimension dependency
162 for multivariate time series forecasting. In *The Eleventh International Conference on Learning*
163 *Representations*.
- 164 Zhou, H., Zhang, S., Peng, J., Zhang, S., Li, J., Xiong, H., and Zhang, W. (2021). Informer:
165 Beyond efficient transformer for long sequence time-series forecasting. In *The Thirty-Fifth AAAI*
166 *Conference on Artificial Intelligence, AAAI 2021, Virtual Conference*, volume 35, pages 11106–
167 11115. AAAI Press.
- 168 Zhou, T., Ma, Z., Wen, Q., Wang, X., Sun, L., and Jin, R. (2022a). Fedformer: Frequency enhanced
169 decomposed transformer for long-term series forecasting.
- 170 Zhou, T., Ma, Z., xue wang, Wen, Q., Sun, L., Yao, T., Yin, W., and Jin, R. (2022b). FiLM: Frequency
171 improved legendre memory model for long-term time series forecasting. In Oh, A. H., Agarwal,
172 A., Belgrave, D., and Cho, K., editors, *Advances in Neural Information Processing Systems*.
- 173