

LLM for time series forecasting

Group: Leonardo Bellini, Lorenzo Grandi, Eugenio Muscinelli

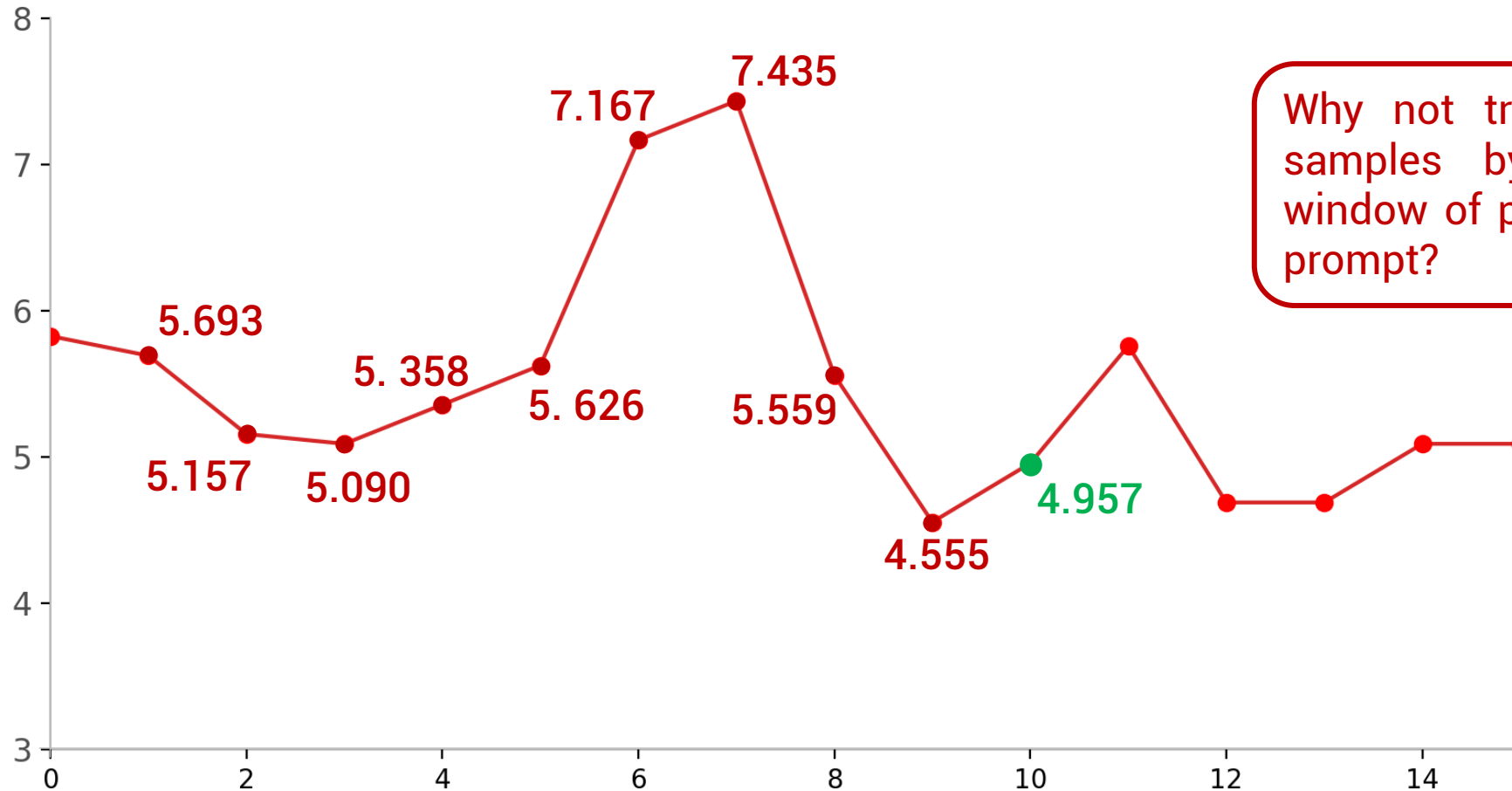
Tutor: Giovanni Battista Esposito

Goal

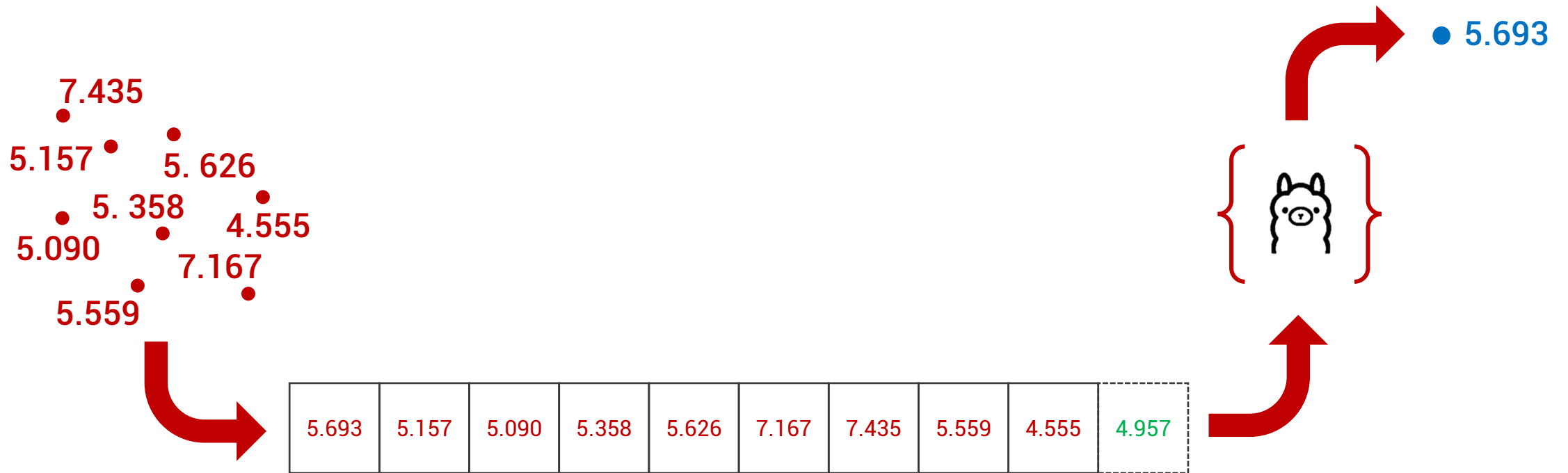
Using the LLama2 model, we want to predict the future samples of a time series representing the data telemetry coming from a node



The starting idea



The starting idea



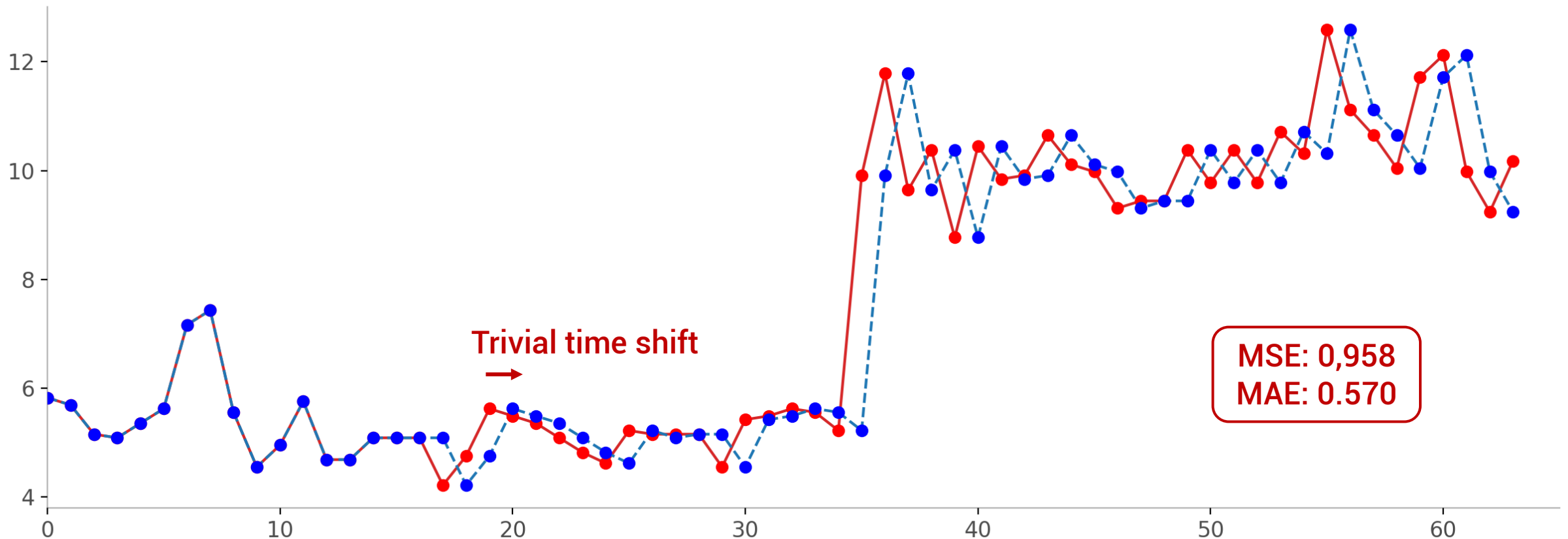
The starting idea

```
1 MODEL_ID = "TinyLlama/TinyLlama-1.1B-intermediate-step-1431k-3T"
2 pipe = pipeline('text-generation', model=MODEL_ID, device=0)
3
4 for batch_idx in range(num_batches):
5     for i in range(batch_idx * batch_size, (batch_idx + 1) * batch_size):
6         window = HUFL_df.iloc[i:i + w_dim]
7         prompt = window[:w_dim - 1].to_string(header=False, index=False)
8         out = pipe(prompt, max_new_tokens=6)[0]['generated_text']
9         pattern = r'\b\d+\.\d{3}\b'
10        out_float = re.findall(pattern, out)
11        out.append(out_float[-1])
```

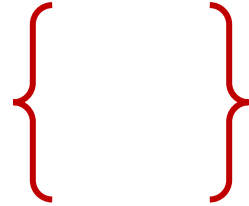
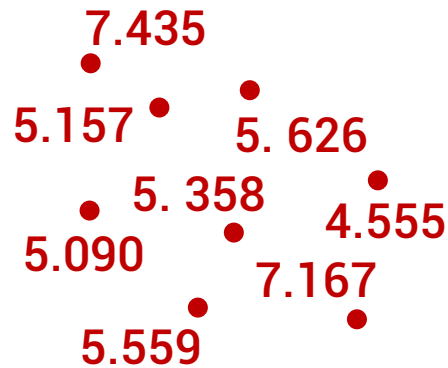
A window of the time-series is passed to the LLM as a prompt

The results

Actual vs Predicted Values



The complete network



Tokenizer : splits the text into individual words (and subwords) and map each token to a unique identifier. In this way, the human-readable text is converted into a form that the LLM can process effectively.

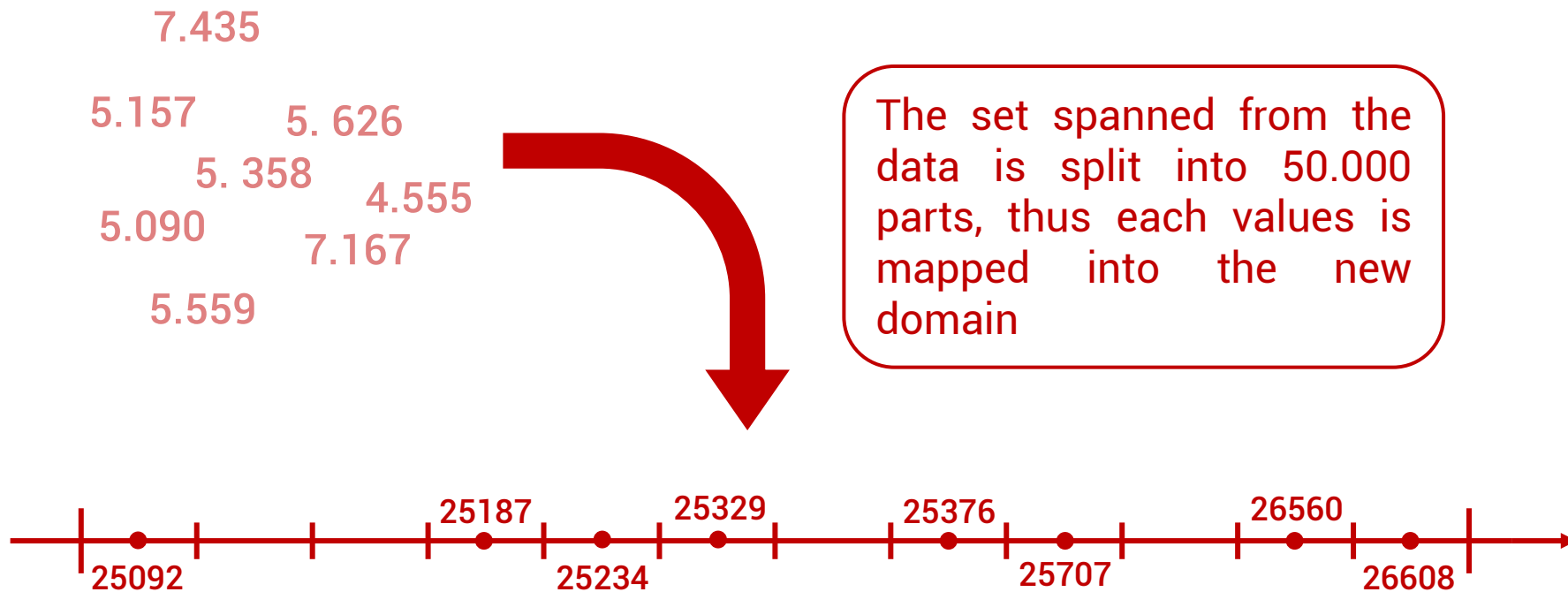
Tokenizer

7.435
5.157 5.626
5.358 4.555
5.090 7.167
5.559

The set spanned from the data is split into 50.000 parts, thus each values is mapped into the new domain



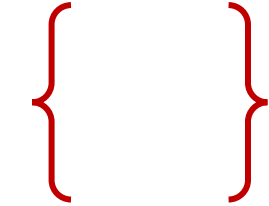
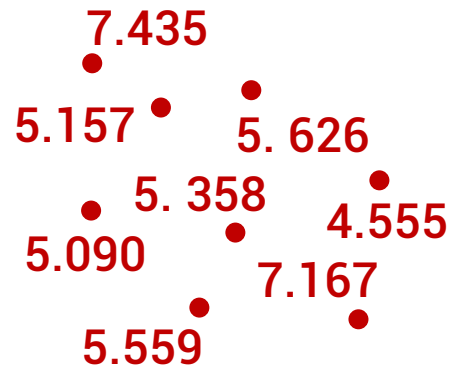
Tokenizer



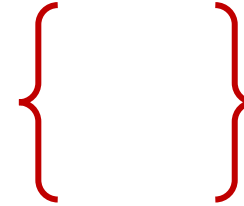
Tokenizer

```
1 dataset = pd.read_csv('/kaggle/input/etth1234/ETTh1.csv')
2 dataset.drop(columns=['date'], inplace=True)
3 scaler = StandardScaler()
4 scaled_df = scaler.fit_transform(dataset) Standardization
5 L = 4000
6 train_df = pd.DataFrame(scaled_df[:L])
7 valid_df = pd.DataFrame(scaled_df[L:2*L])
8 test_df = pd.DataFrame(scaled_df[2*L:3*L])
9 train_HUFL = torch.Tensor(train_df[0].values)
10 valid_HUFL = torch.Tensor(valid_df[0].values)
11 test_HUFL = torch.Tensor(test_df[0].values)
12
13 class Tokenizer(nn.Module):
14     def __init__(self, voc_size, lower_bound, upper_bound):
15         super(Tokenizer, self).__init__()
16         self.voc_size = voc_size
17         self.lower_bound = lower_bound
18         self.upper_bound = upper_bound
19         self.discrete_values = torch.linspace(self.lower_bound, self.upper_bound, self.voc_size).to('cuda')
20
21     def forward(self, dataset_std):
22         quantized_values = (torch.bucketize(dataset_std, self.discrete_values) + 2) Quantization
23         attention_mask = (~torch.isnan(dataset_std))
24         return quantized_values, attention_mask
```

The complete network

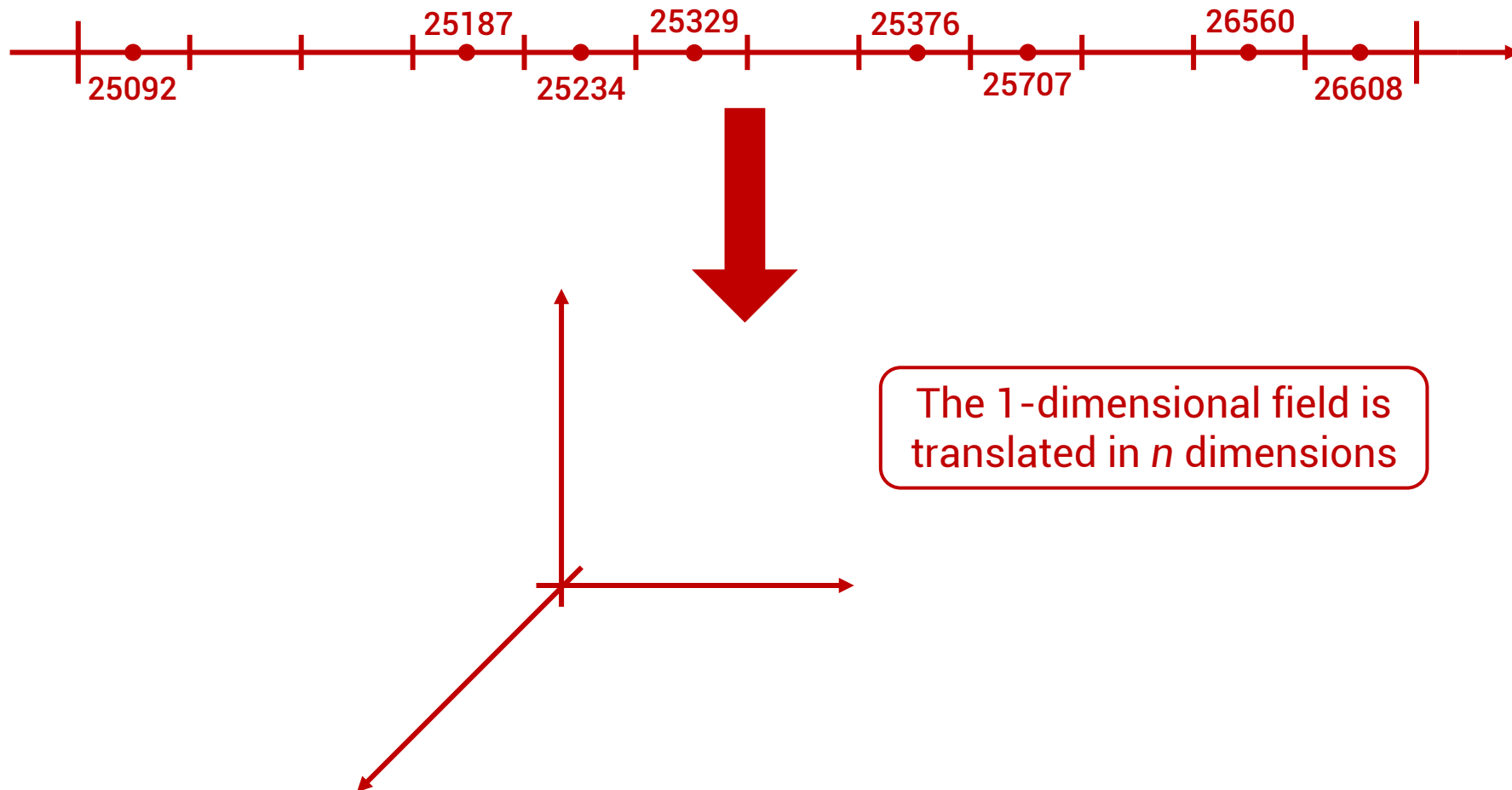


Tokenizer

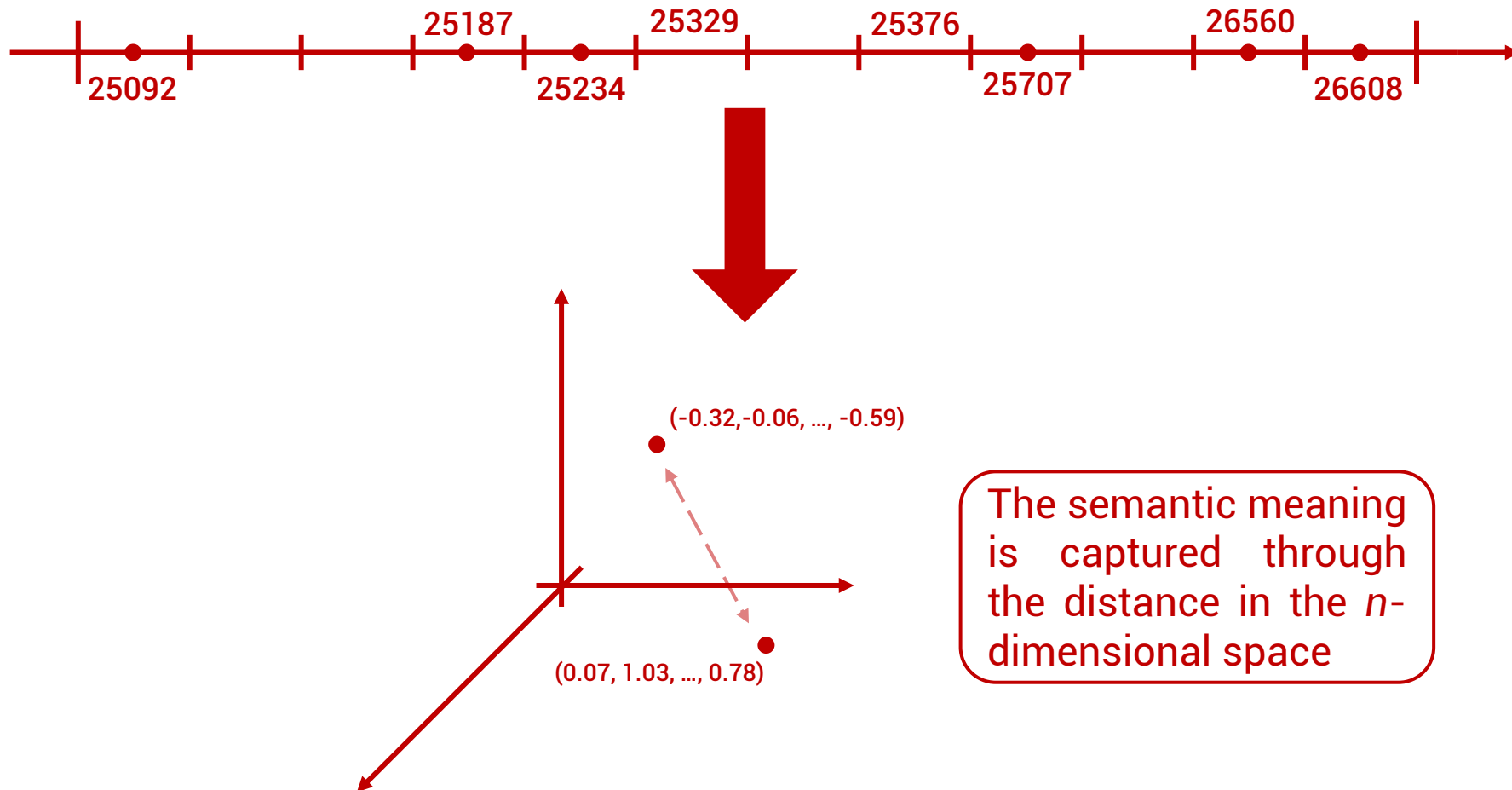


Word embedder : transforms the tokens into dense vectors of fixed size, named embeddings, which capture the semantic information about the words

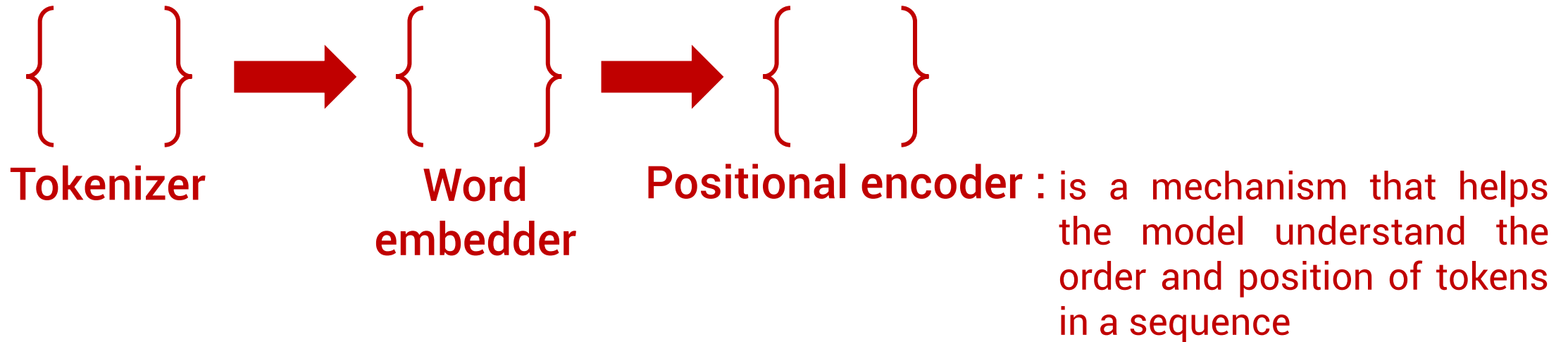
Word embedder



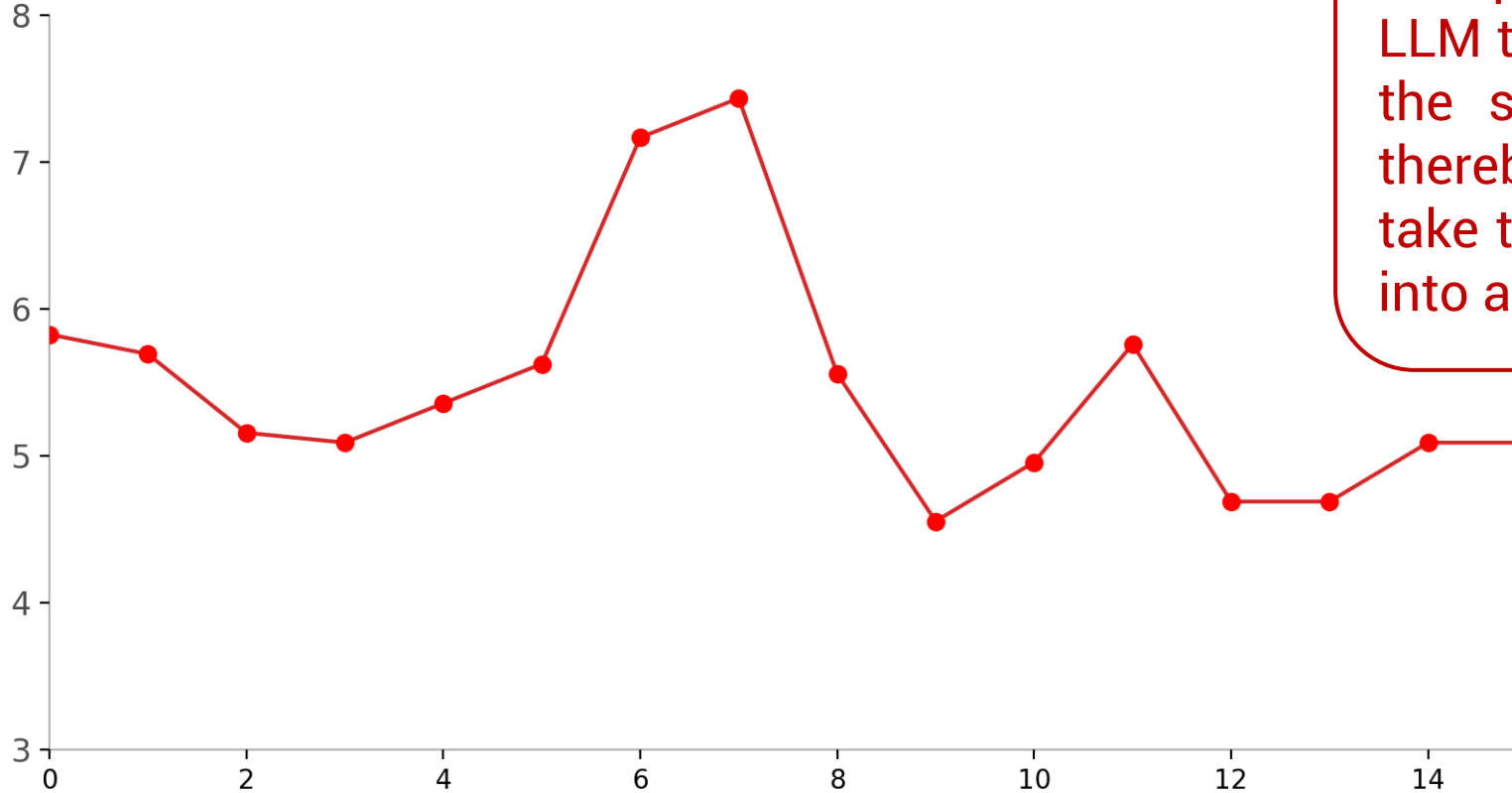
Word embedder



The complete network

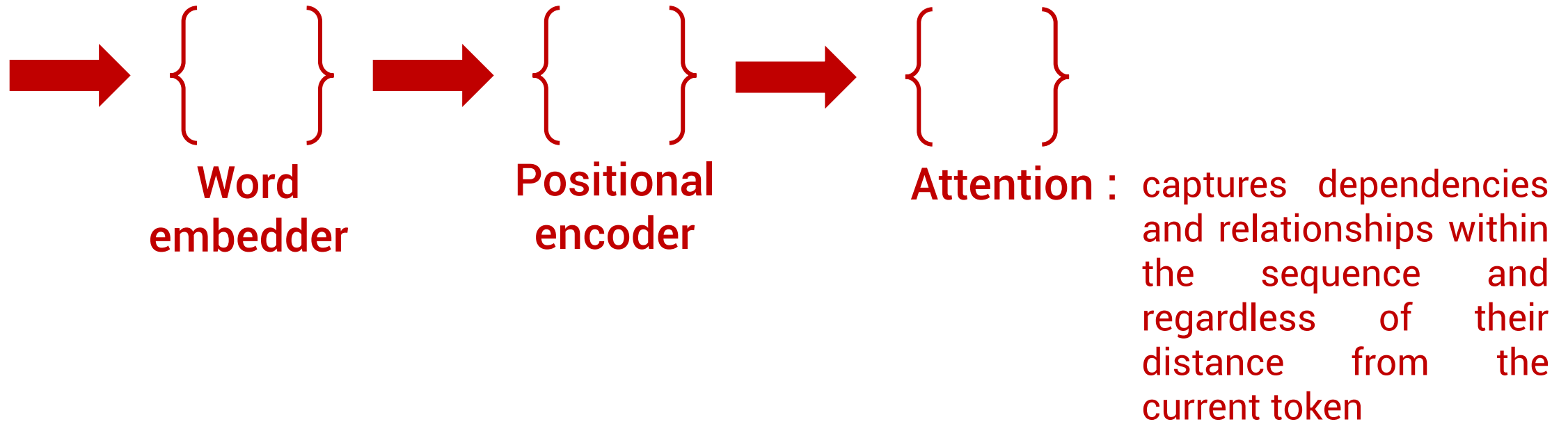


Positional encoder

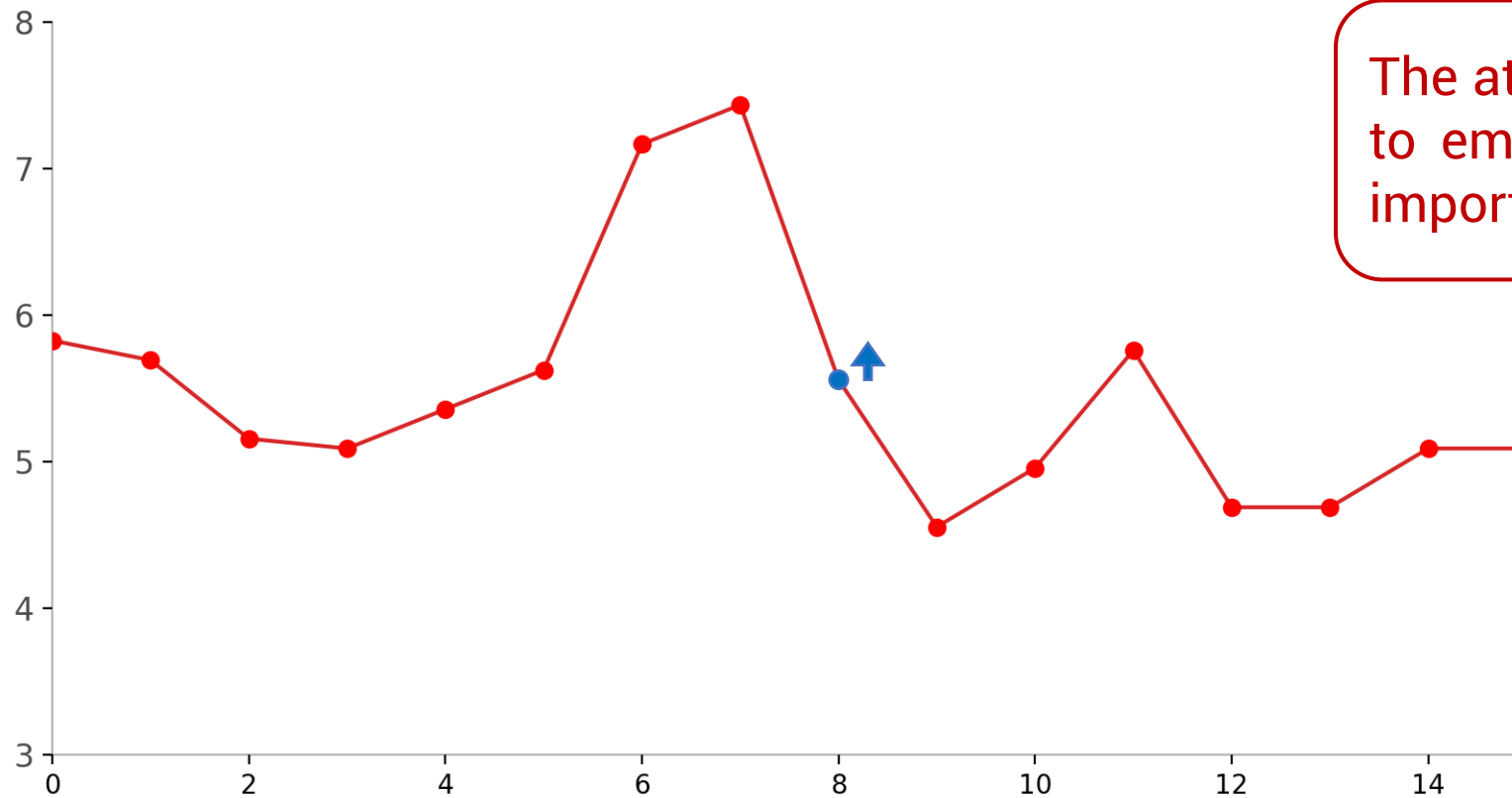


The positional encoder allow the LLM to keep the semantic order of the samples in the time-series, thereby enabling predictions that take the history of the time-series into account

The complete network

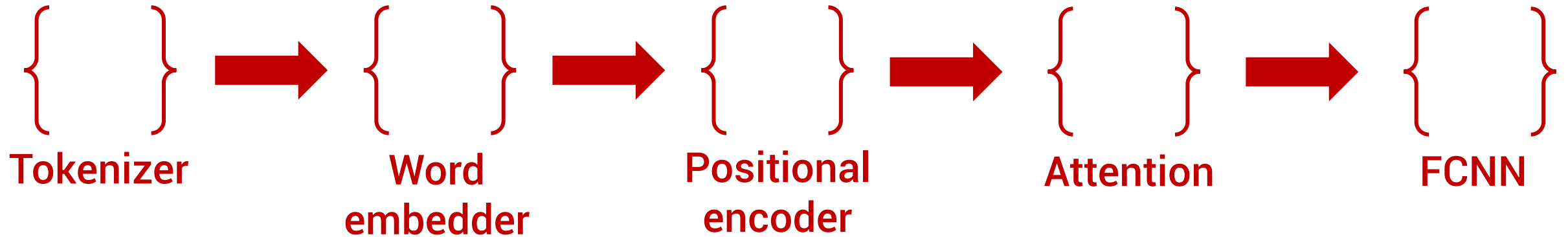


Attention layer



The attention layer assigns a score to embedding token related to its importance in the context

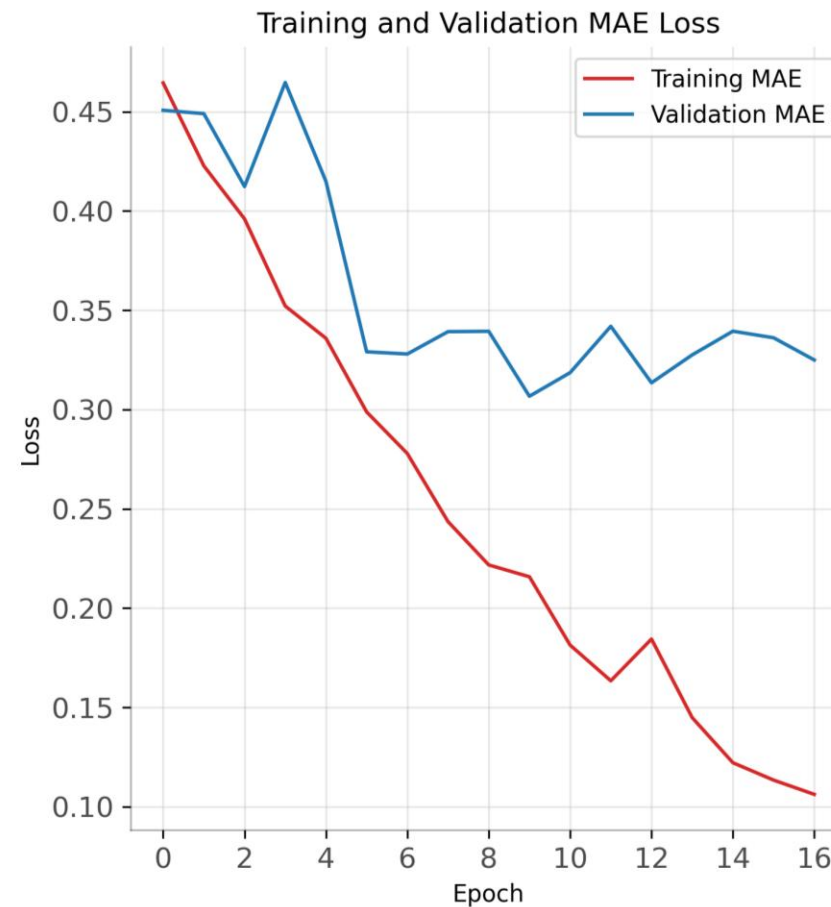
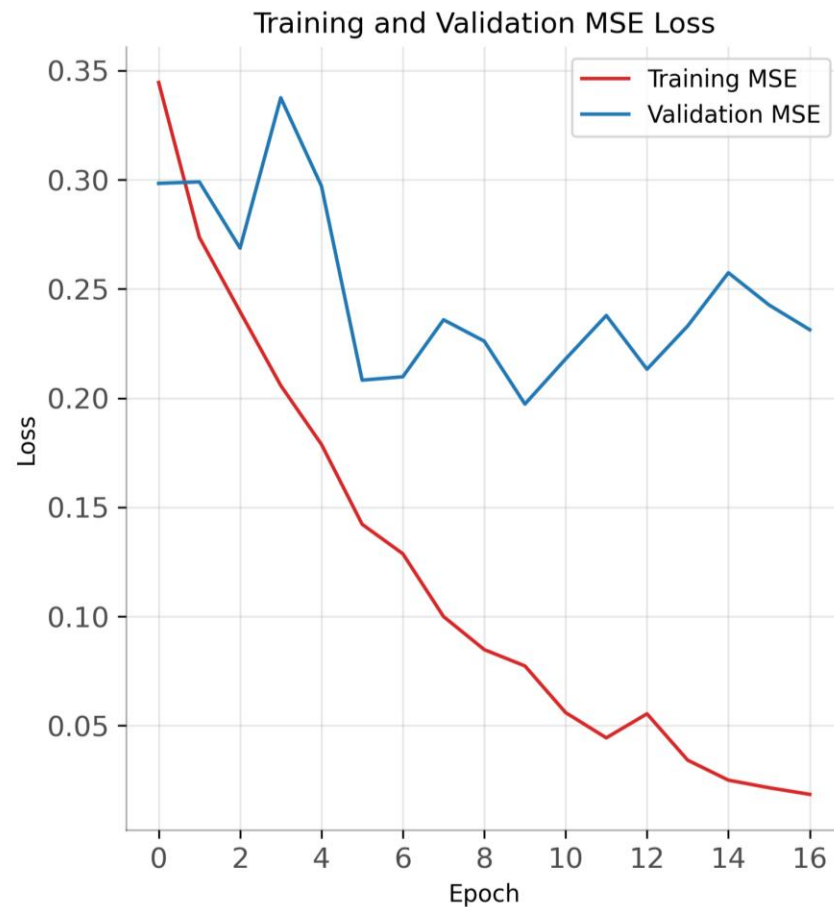
The complete network



Two Fully-Connected NN are trained to map the attention embedding to the predicted value, acting as a detokenizer

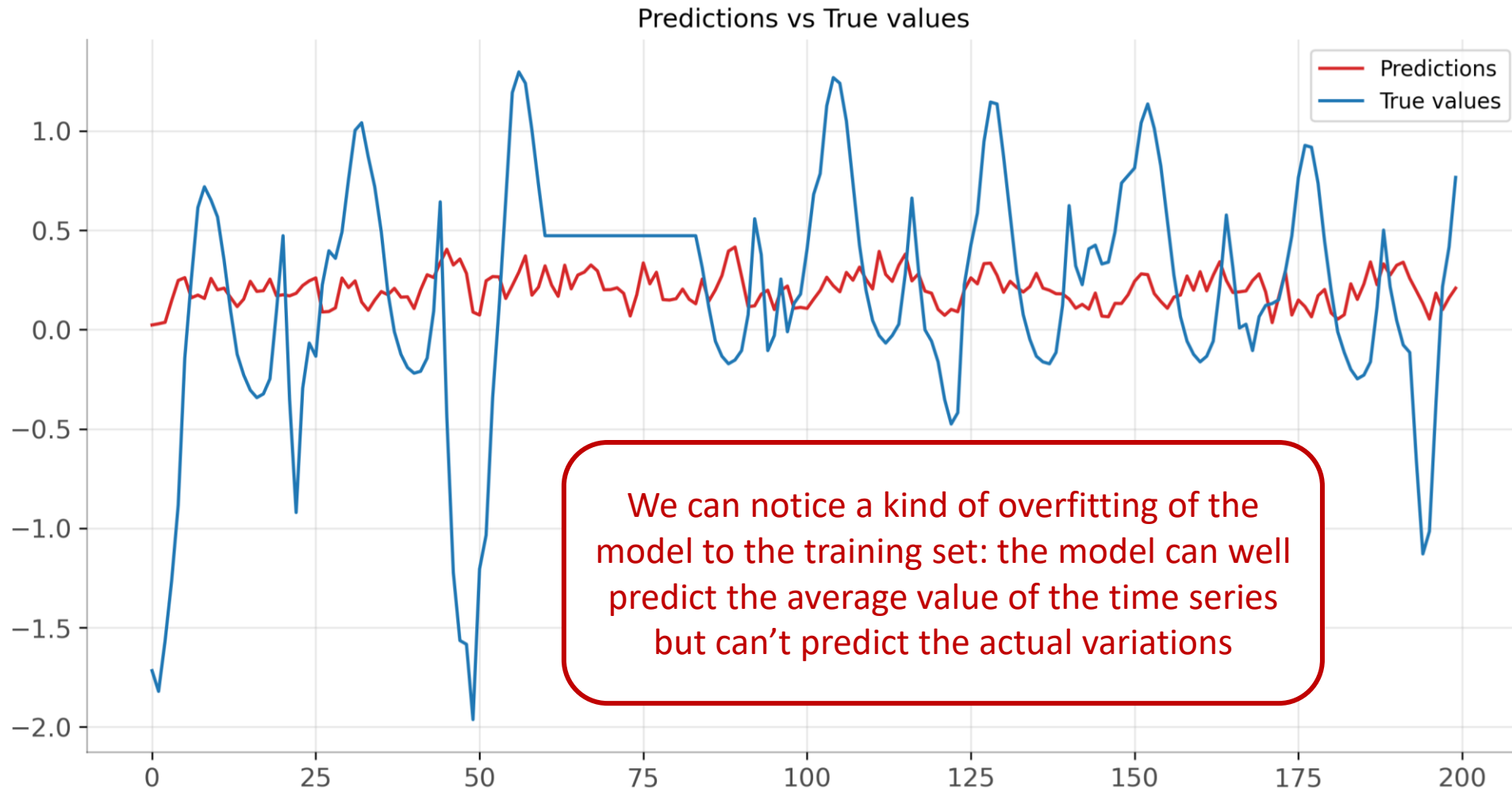
The results

NB: Tested on a second dataset,
uncorrelated to the one used in the
training phase

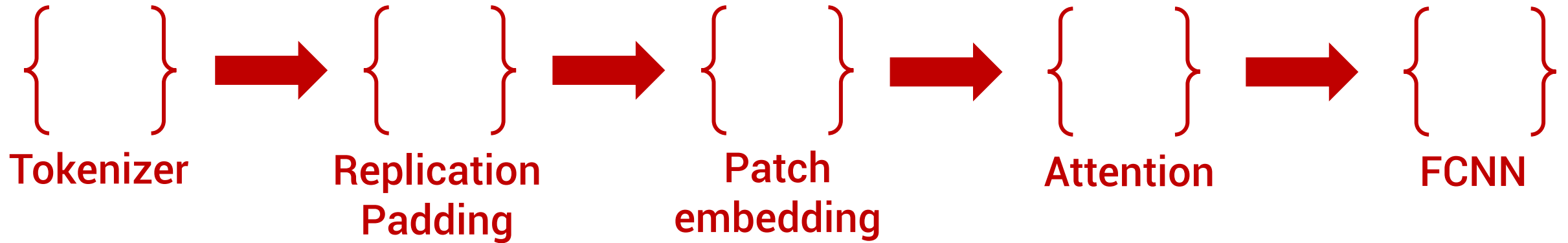


MSE: 0,524
MAE: 0,466

The results

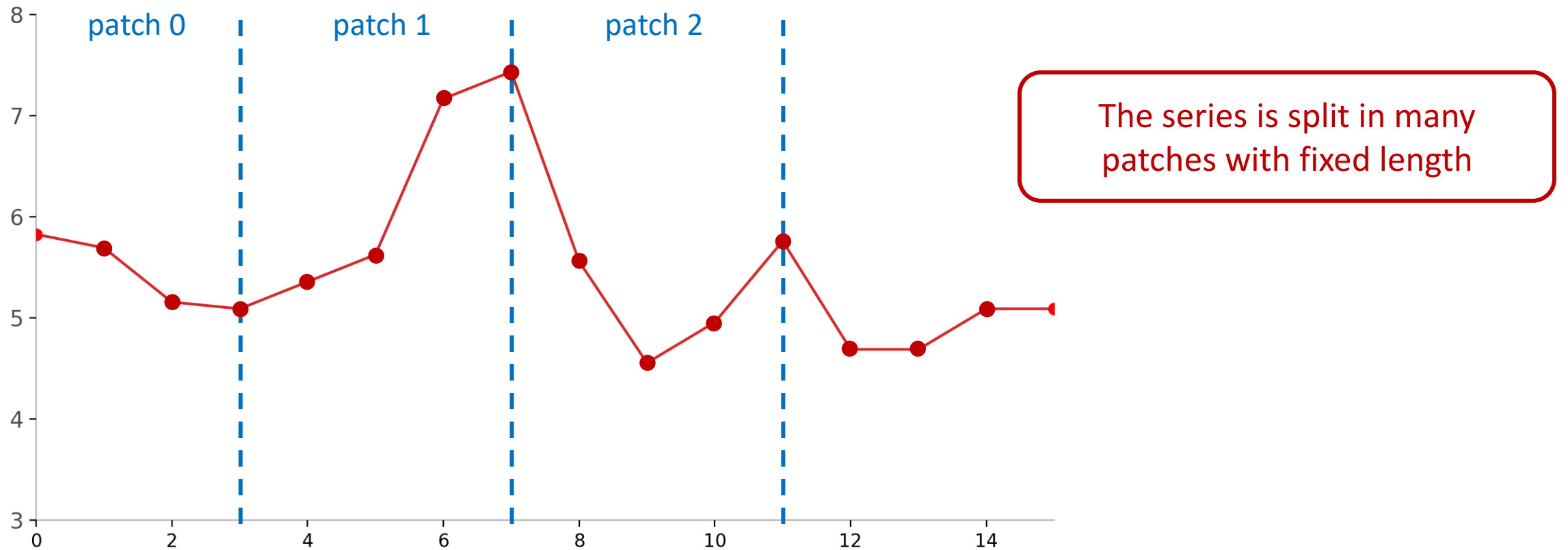


The alternative solution



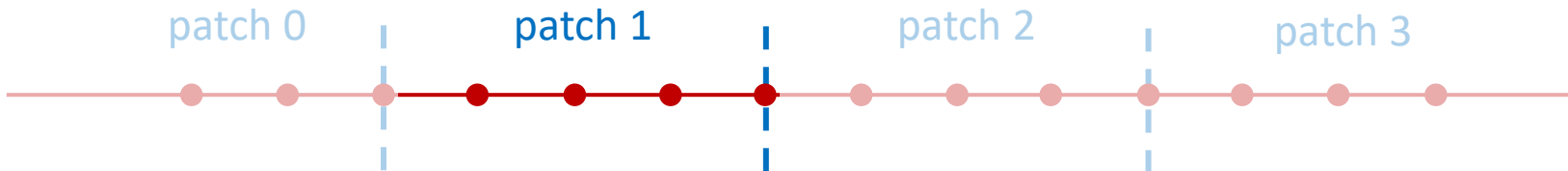
(Ref: Time-llm: time series forecasting by reprogramming large language models)

Patch embedding



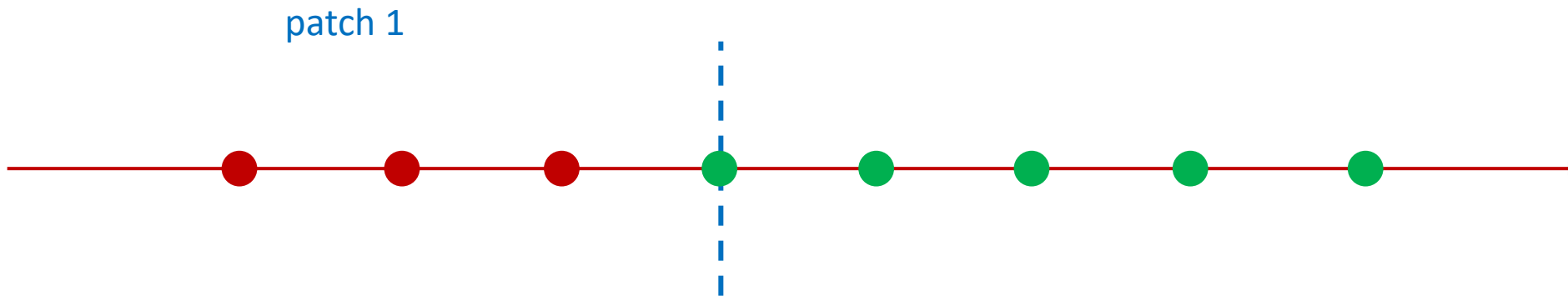
Patch embedding

The series is split in many patches with fixed length

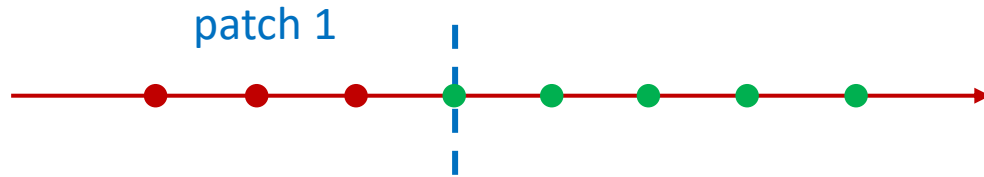


Patch embedding

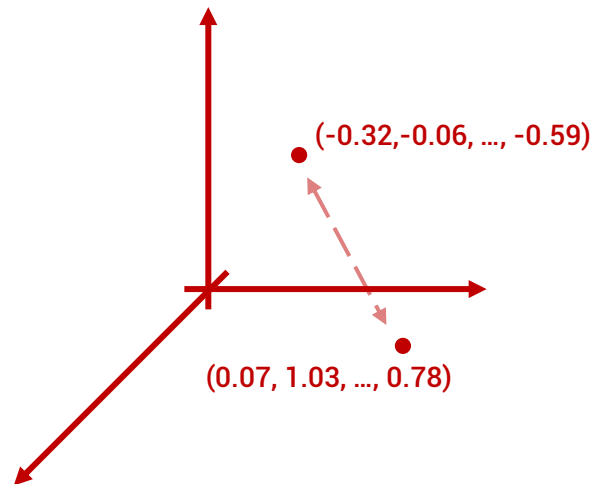
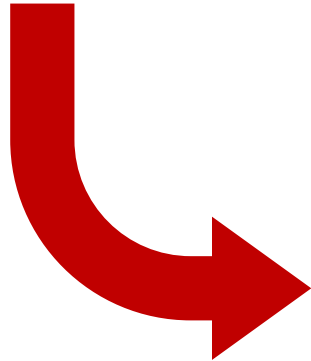
The previous layer pad
and stride the patch,
adding the last value to
highlight it



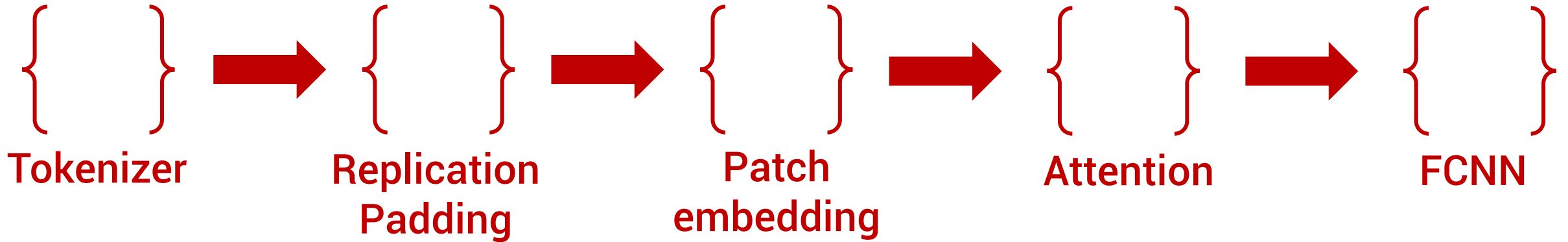
Patch embedding



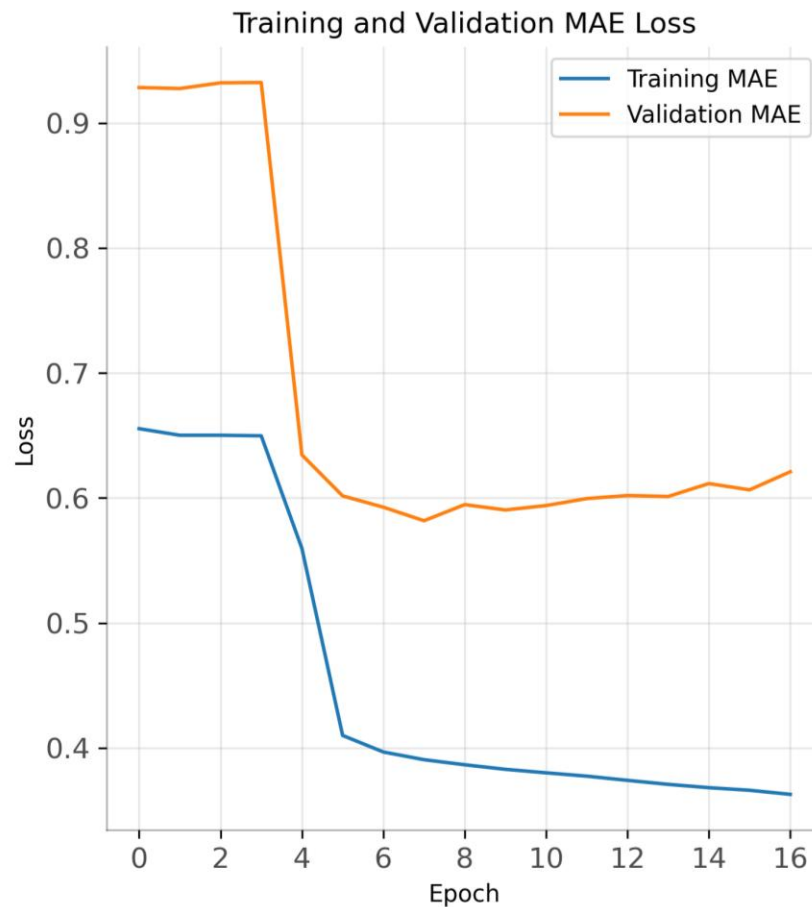
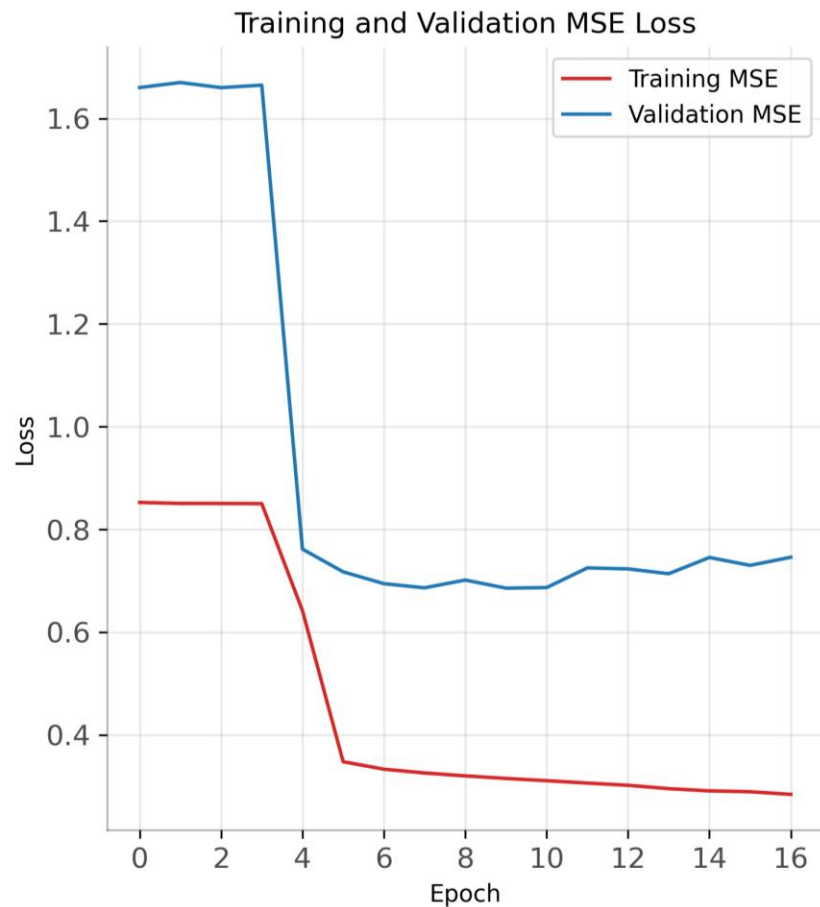
The embedding vectors are computed through a monodimensional convolution



The alternative solution

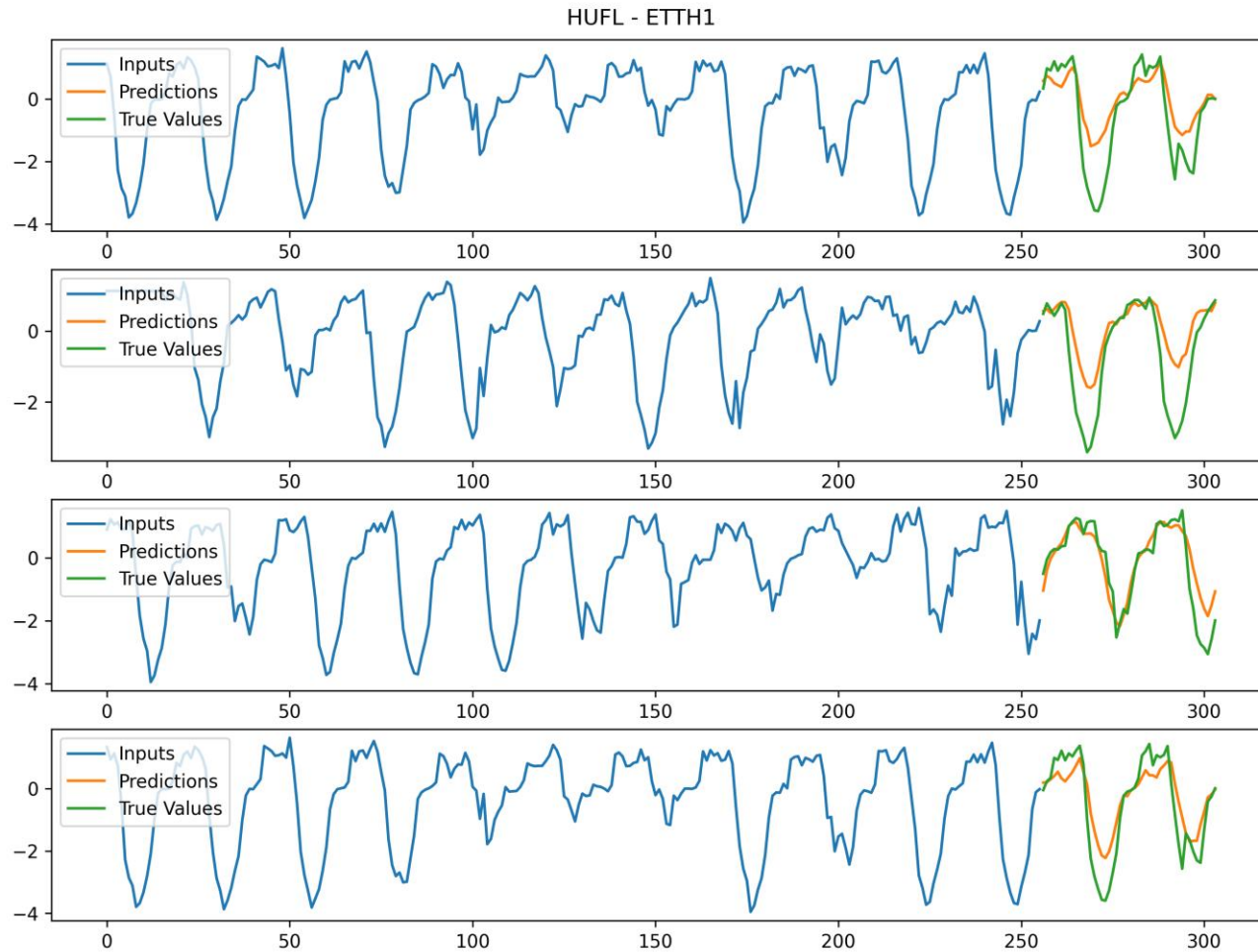


The results



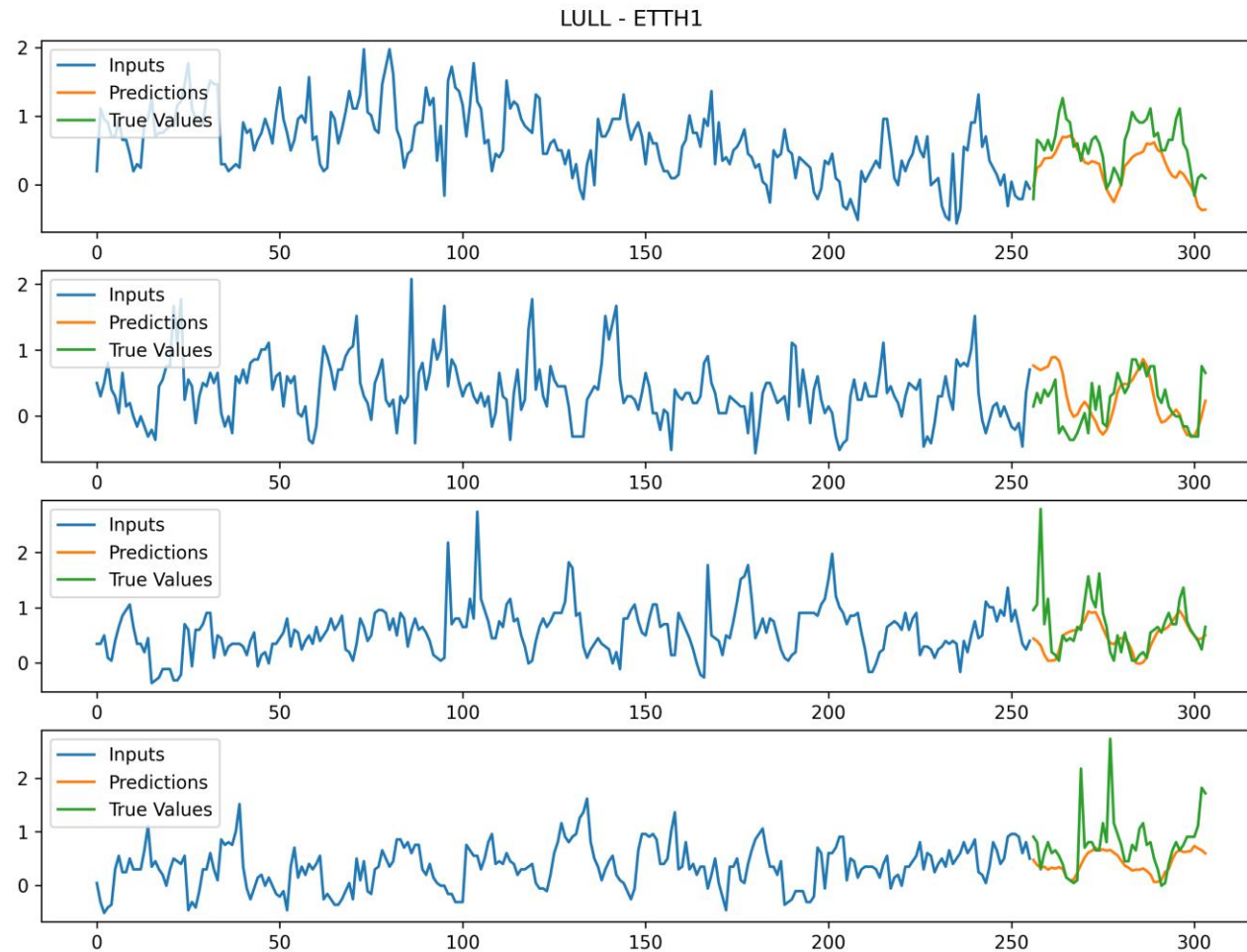
MSE: 0,536
MAE: 0,522

The results



Even with the smallest
model the prediction fit
the true behavior

The results



Using models with a larger number of parameters would result into better predictions

The results



We tried to evaluate the performances on a different dataset, and the result isn't that different with respect to the previous tests

Possible improvements

From the model viewpoint:

- Larger pretrained model
- Quantization of the parameters
- More complex layers

From the system viewpoint:

- Split the network processing over more GPUs

Thank you for your attention
