

# Prefix Tuning

## Instructor

Sourab Mangulkar

Machine Learning Engineer at   
Creator of  PEFT



# Prefix Tuning

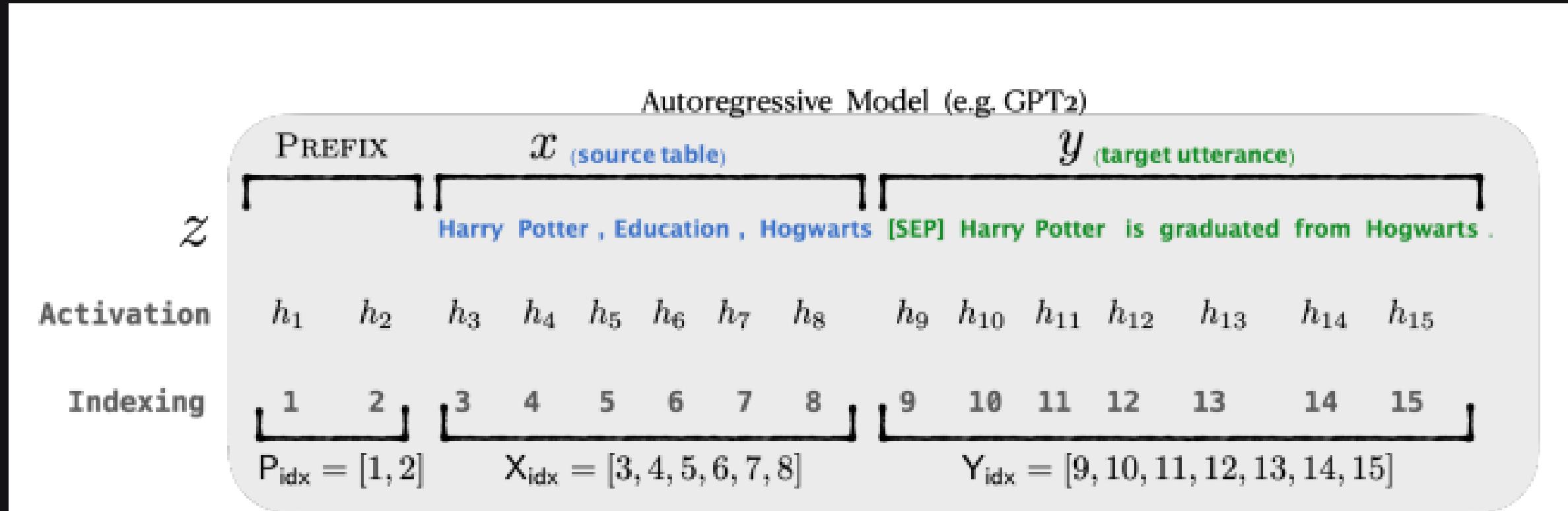


FIGURE EXPLAINING HOW PREFIX-TUNING WORKS FROM [ORIGINAL PAPER](#), FIGURE 2

# Prefix Tuning

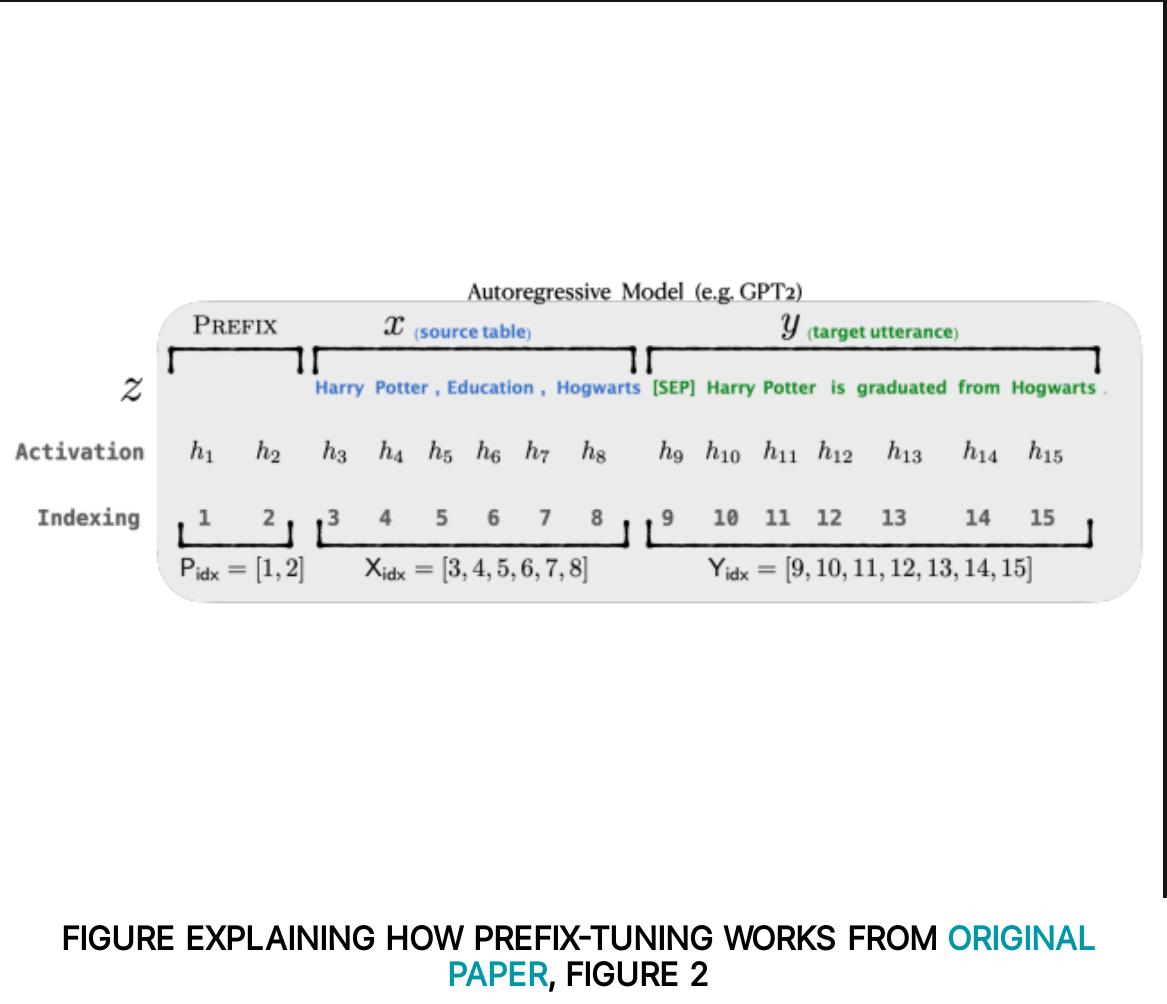


FIGURE EXPLAINING HOW PREFIX-TUNING WORKS FROM ORIGINAL PAPER, FIGURE 2

```
class PrefixTuningAttentionLayer(torch.nn.Module):
    def __init__(self, base_attention_layer, num_prompt_tokens, emb_dim, batch_size):
        super().__init__()
        self.base_attention_layer = base_attention_layer
        self.prompt_embedding = torch.nn.Embedding(num_prompt_tokens, emb_dim)
        self.prompt_tokens = (torch.arange(num_prompt_tokens)
                             .unsqueeze(0)
                             .repeat(batch_size, 1)
                             .long())

    def forward(self, hidden_states, **kwargs):
        # get the soft prompt embeddings
        soft_prompts = self.prompt_embedding(self.prompt_tokens)
        # prepend the soft prompt embeddings to hidden states
        hidden_states = torch.cat((soft_prompts, hidden_states), dim=1)
        return self.base_attention_layer(hidden_states=hidden_states, **kwargs)
```