# Transformer-based Seq2Seq and Sparsity & Linearity
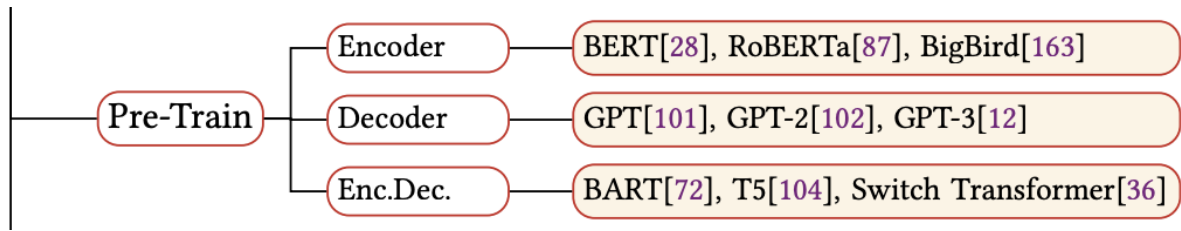
| | |
|---|---|
| ☰ Author | @Dyan Lee |
| 👤 Created By | Ⓓ Dyan Lee |
| 🗓 Date | @2021년 9월 10일 |
| ☰ Organization | KAIST EdLab. |
| ◔ Presented @ | |
| ◔ Published | Sep 2021 |
| ◔ Tags | Kick-Off |
| ☰ Tags 1 | Study |
| ☰ Transcript | |
| ☰ Type | Weekly Sync |

# Agenda

- Walkthrough a survey papers on Transformer, Efficient Transformer, Visual Transformer, and etc.

- Seek out what must be done first.

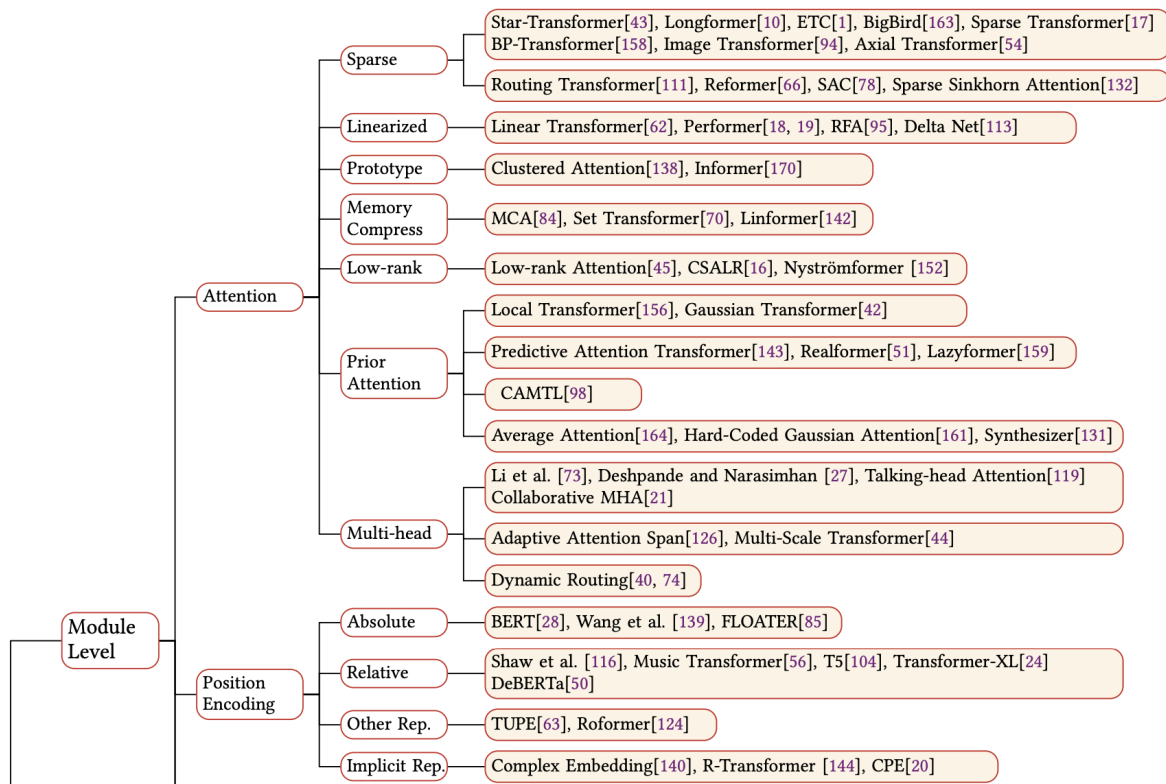- Primary agenda will be Performer and more.

# Summary

## Seq2Seq Models

## Pre-Train

- BART[72], T5[104], Switch Transformer[36]

- **gpt**

# Attention



# Sparsity

- **Star-Transformer**[43], **Longformer**[10], **ETC**[1], **BigBird**[163], **Sparse Transformer**[17]
  BP-Transformer[158], **Image Transforme**r[94], Axial Transformer[54]

- Routing Transformer[111], **Reformer**[66], SAC[78], **Sparse Sinkhorn Attention**[132]

## Linearised

- **Linear Transformer[62], Performer[18, 19], RFA[95], Delta Net[113]**


1. **Linear Transformers**

   - Feature Maps에 따른 분류.

   - Feature space 상에서의 orthogonality를 이용할 수 있는 feature map을 제안했다.

2. Performers

   - 

## Prototype

- Clustered Attention[138], **Informer**[170]

## Memory Compress

- MCA[84], Set Transformer[70], **Linformer**[142]

# Multi-model