# Transformer□□

1. □□□□□□□□□□□□
2. □□□□□□□□□
3. scaling□□□□□□□□□□□□□□□□□□□□□□□
4. "□□□□ "□□□□□□ GPU□□□□□□ GPU□□□
5. □ nlp□□□□□□□□□□□□□□ token□ □□ □□□□□
6. □□□□□□□□□□□□□ token□□□□
   1. □□ kq□□□□□□□ token□□□□□□□□ v□□□□□□□
   2. □□□□□□□□ kqv□□□□□□□ sequence□□□
   3. □□□□□□□□□□□□□□□□□
   4. weight□□□□□□□□□□□□□□□□□□□□□□

## □□

1. □□□□□□□□□□□□□□□□□□
2. □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
   □□□□□□ AGI□□□□□□□□ AGI□□□□□□□□□□□□□□□□□□□□
3. Transformer
   □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

# Attention

1. □□□□ token□□□□□□□□□□ token□ key□ value
   1. □□□□□□□□□□□□□□ □□ □□□
      1. □□□□□□□□□□□□□□□□□□□□□□□□
      2. □□□□□□□□□□□□□□□□□□ □□□□□ □□□
   2. □□□□□□ KV cache□□□□□□□
2. □□□□□□□□□
3. □□□□
4. □□□□□□□□□□ attention□□
   1. □□□□□□□□□□□□□□□□□□□□□□□ O(T)□ O(1)□□
      1. Transformer □□□□□□□ O (T^2)□□□□□□ O (T^2)
      2. RWKV□□□□□ status□□□□□□□□□□□□□□□□ ,□□□□□□
         O(T)□□□□□ O(1)
   2. □□□□□□□ token
      □□□□□□□□□□□□□□□□□□□□□□□□□
   3. □□□□□□□□
      1. MOE□ CoT□□□ □□□□□□

5. □□□□□□□□ FFN□□□□□
6.

# System2

1. Transformer□□□□□□ System 2□□□ RL+CoT□□ □□□□□□
2.

# □□ API

1. transformer□□□□□□□□□ attention□□□□□ □□□ □□□□□□□□□□□□
2. attention□□□□□□□□□□□□□□□□□□□□□□□
3. □□□□□□□□□□□□□□ token□□□□□□□□□□□□□□□□
4. □□□□□□□□□□□□□□□□□□
5. □□□□□□□□□□□□□□□□□□
6. □□□□
   1. □□□□□□□□□□ transformer□□□□□□□□□
   2. □□□□□□□□□□□□□□□□□□□□□□□□□□□□

# □□□□

LLM
□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
□□□□□□□□□□□□□□□□□□□□□□□□□□

□□□ LLM
□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
□□□□□□□□□□□□□□□□□□□□□ □

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
□ AI□□□□□□□□□□□□□□□□□□□

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
□□□□□□□□□□□□□□□