

# Transformer

??

1. [ ]
2. [ ]
3. scaling [ ]
4. " [ ] " [ ] GPU [ ] GPU [ ]
5. [ ] nlp [ ] token [ ] [ ] [ ]
6. [ ] token [ ]
  1. [ ] kq [ ] token [ ] v [ ]
  2. [ ] kqv [ ] sequence [ ]
  3. [ ]
  4. weight [ ]
7. [ ]
  1. LLM [ ]
  2. LLM [ ] [ ] [ ]
8. hidden size 7168 [ ]
  1. [ ]
  2. [ ]
  3. 7168 hidden [ ] token [ ]
    1. [ ] hidden [ ] 1024?) [ ] KVcache [ ]
    2. [ ] sequence [ ] token [ ]
9. [ ]
  1. [ ]
  2. [ ]
10. [ ]
  1. [ ]
  2. [ ]

??

1. [ ]
2. [ ] [ ] AGI [ ] AGI [ ]
3. Transformer [ ]

## Attention

1. token key value
  1. KV cache
  1. RWKV
  2. KV cache
2. RWKV
3. RWKV
4. attention
  1. Transformer  $O(T^2)$  RWKV  $O(T)$ 
    1. Transformer  $O(T^2)$  RWKV  $O(T)$
    2. RWKV status  $O(1)$
  2. token
  3. MOE CoT
5. FFN
6. " " " "
  1. attention + MLP
  2. Yan2.0 Preview RockAI Transformer
  1. RockAI
  2. RockAI
  3. kv cache kvcache
  4. kv cache kvcache

## ??API

1. transformer attention
  1. Token CoT Latent Chain-of-Thought
  2. attention
  3. token
  - 4.
  - 5.
  6. transformer

## ?????

