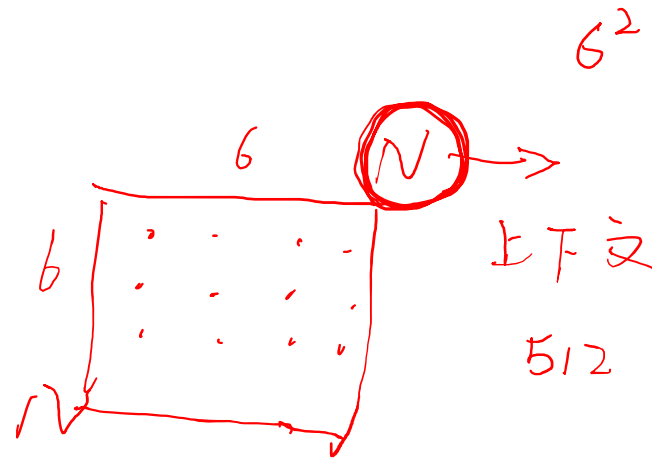


Transformer架构优化

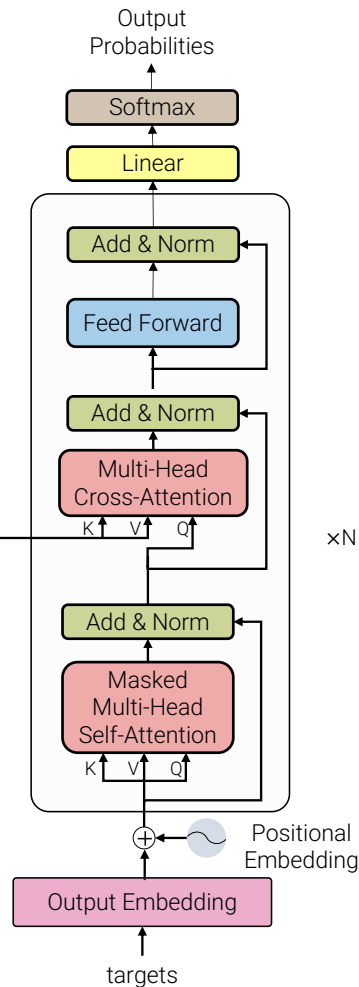
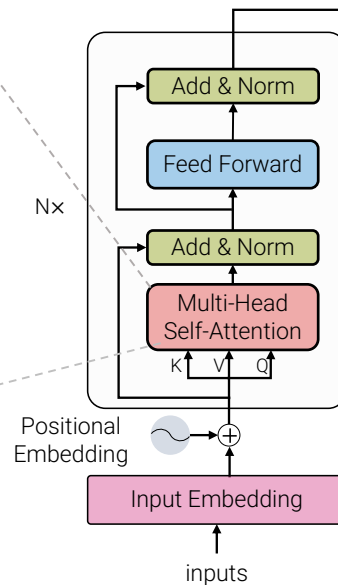
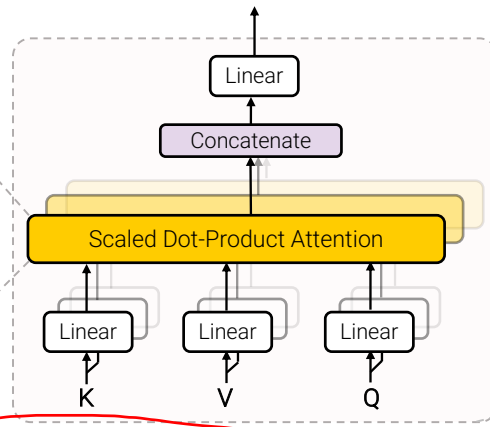
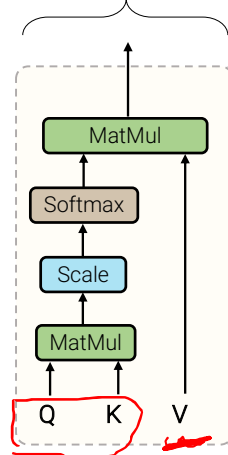
Scott



计算复杂度



Computational and Memory Complexity
 $O(n^2)$

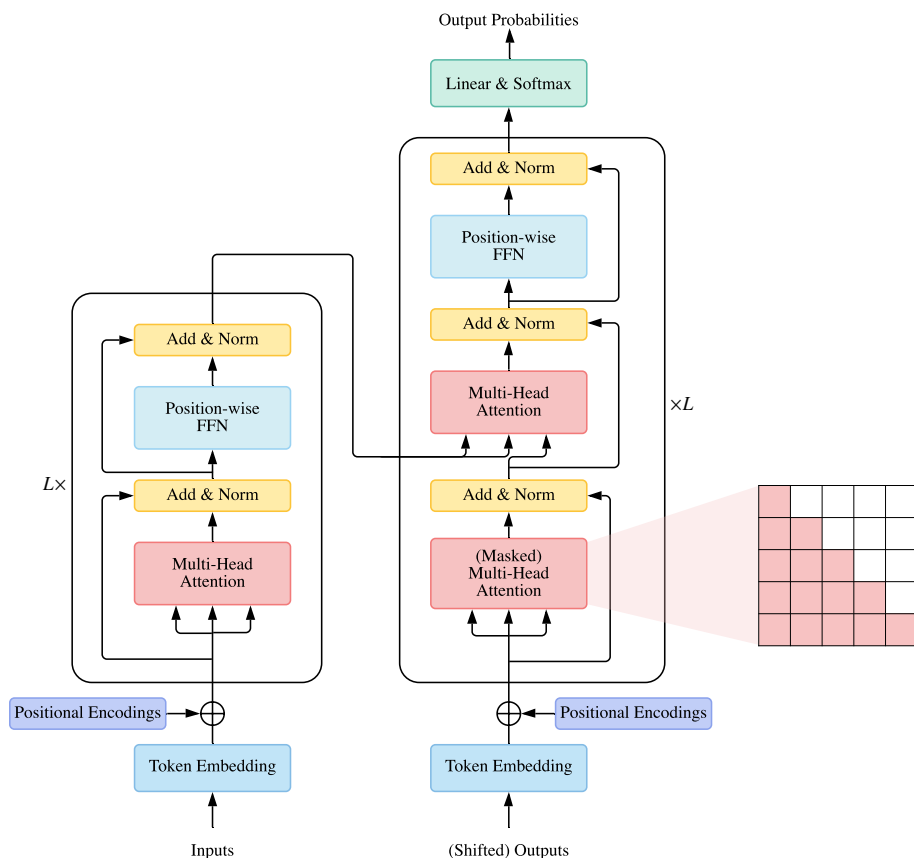


Handwritten notes and calculations:

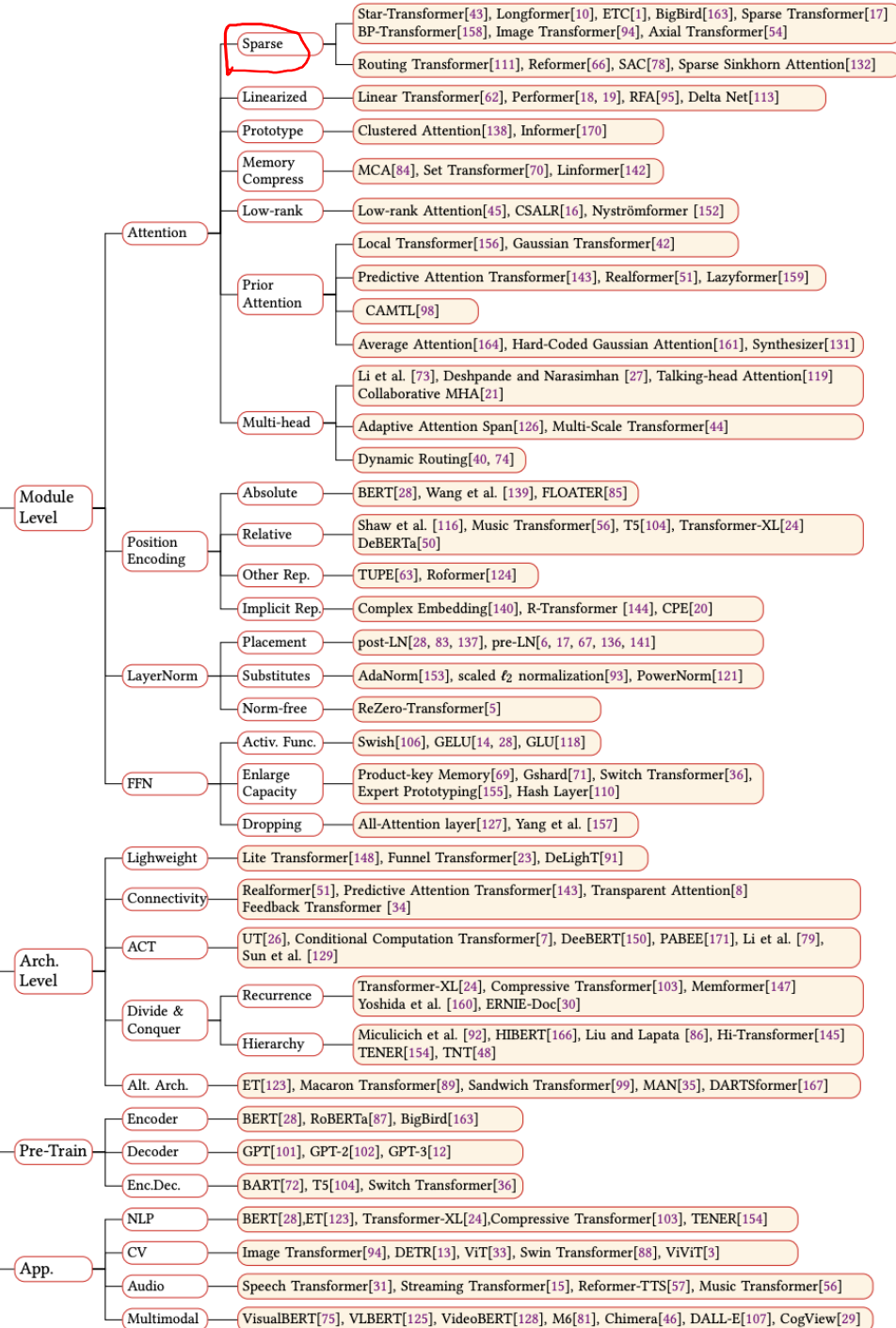
- Q and K^T are 6×512 and 512×6 respectively.
- V is 6×512 .
- The calculation $6 \times 6 \times 512$ is circled in red.
- The text "6x512" is written below the circled calculation.



Transformer家族 → x-former



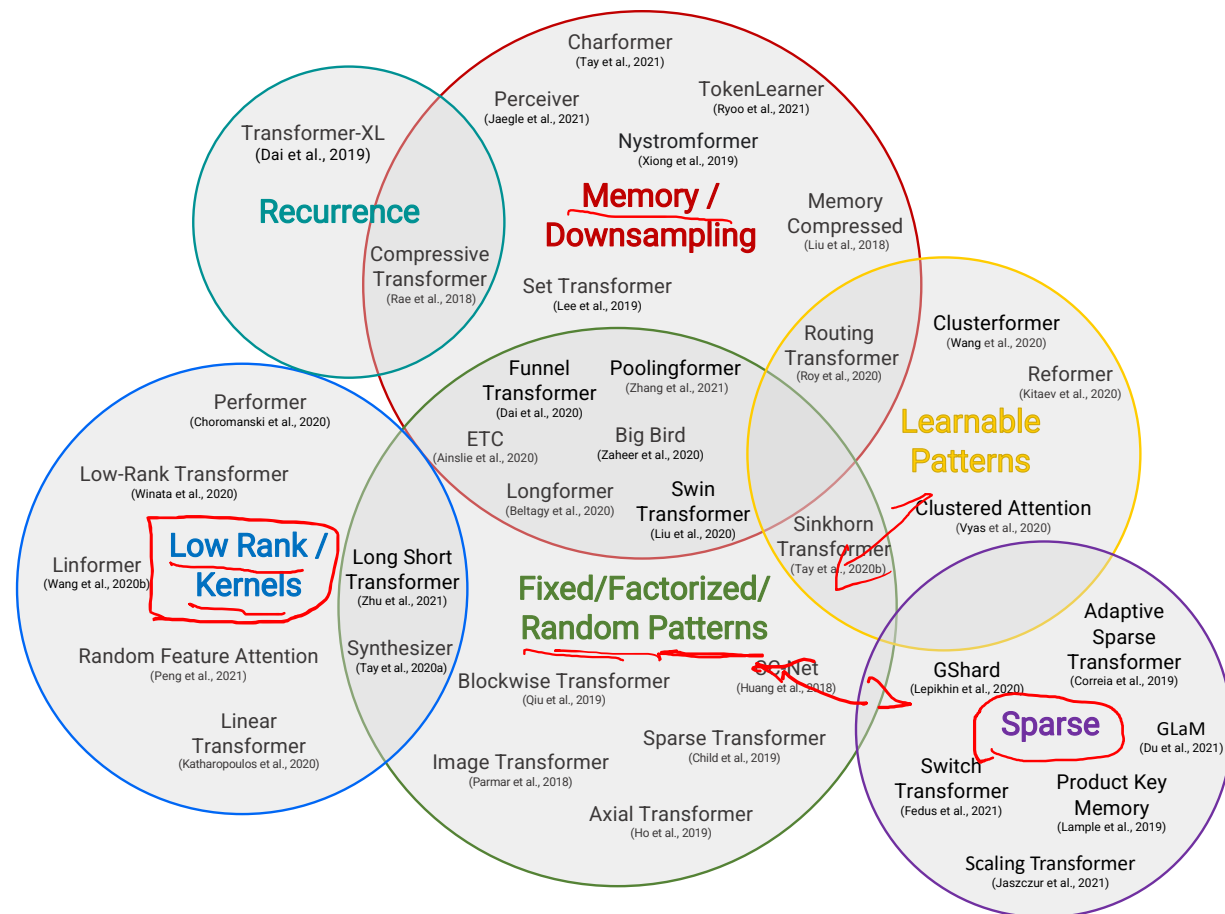
X-formers



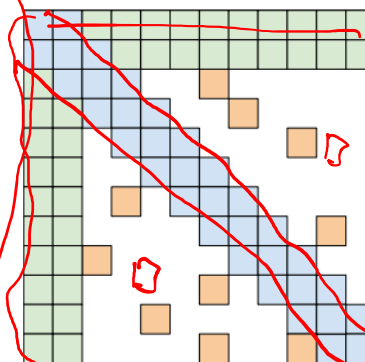
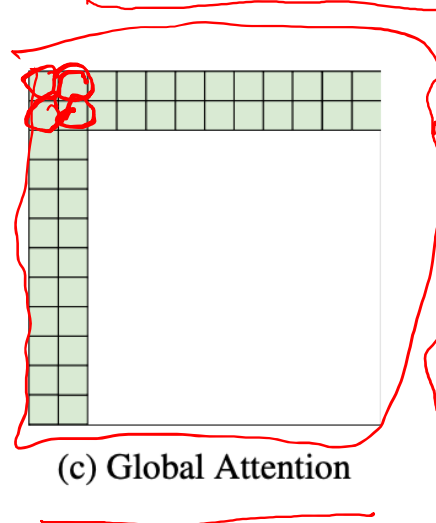
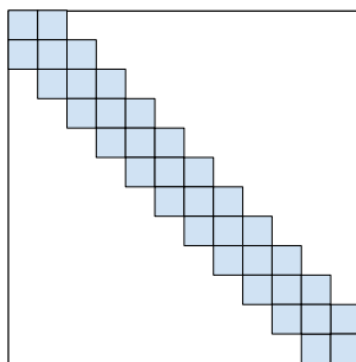
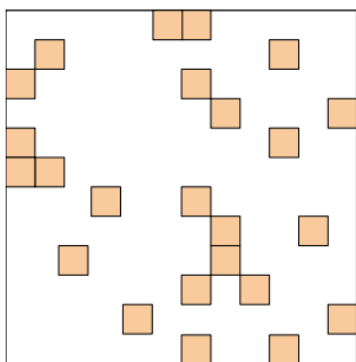
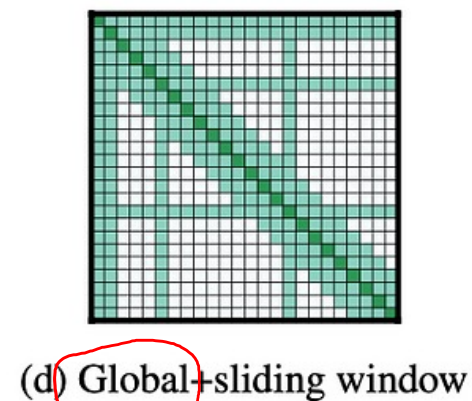
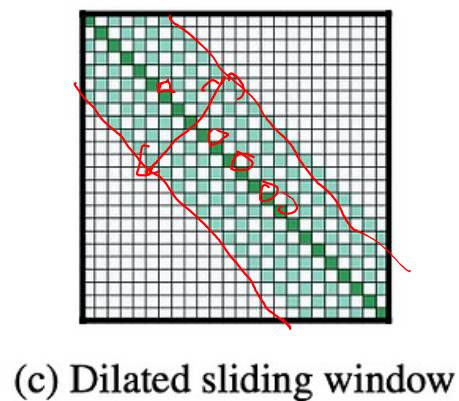
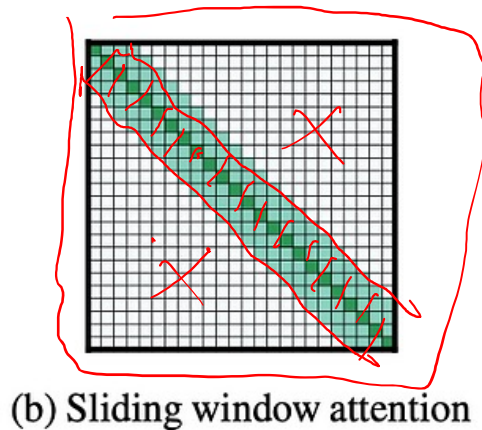
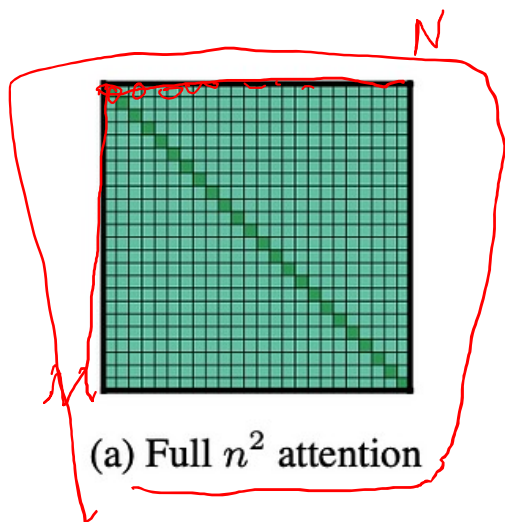
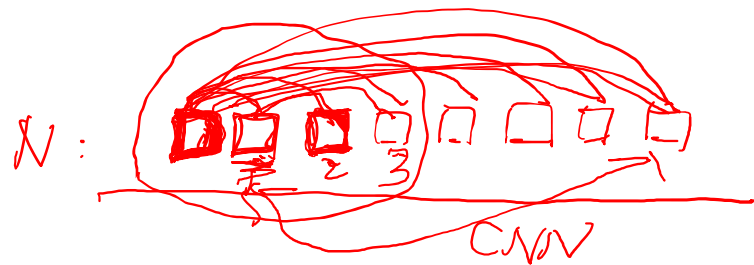
高效Transformer架构设计

Efficient Transformer

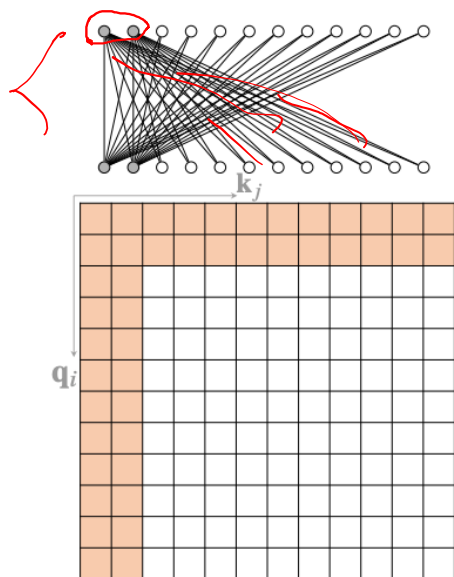
- 增加上下文长度
- 提高速度
- 减少内存使用



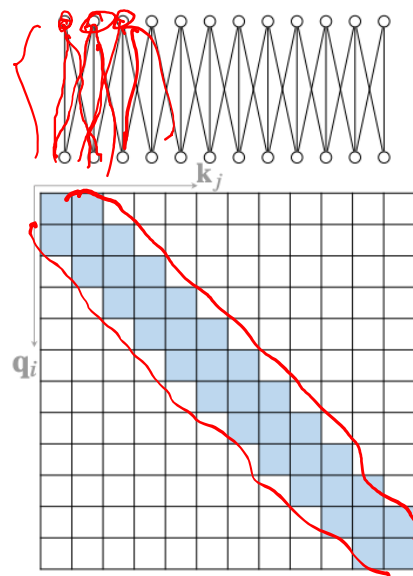
稀疏方法



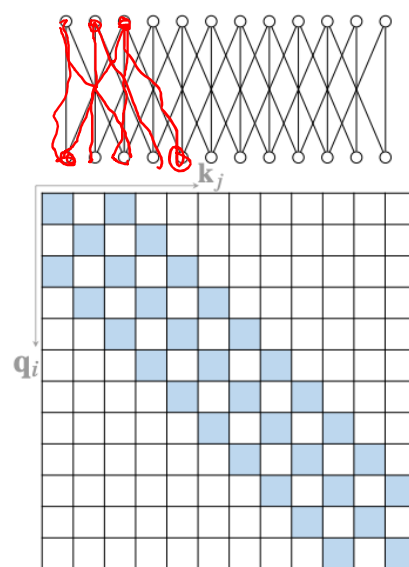
稀疏方法



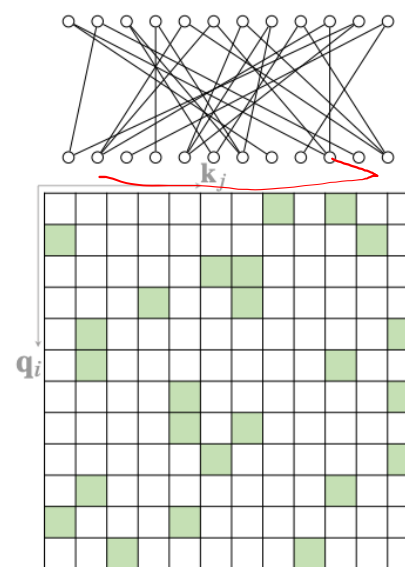
(a) global



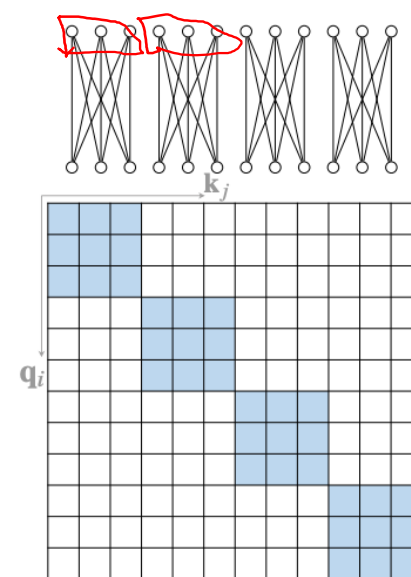
(b) band



(c) dilated



(d) random

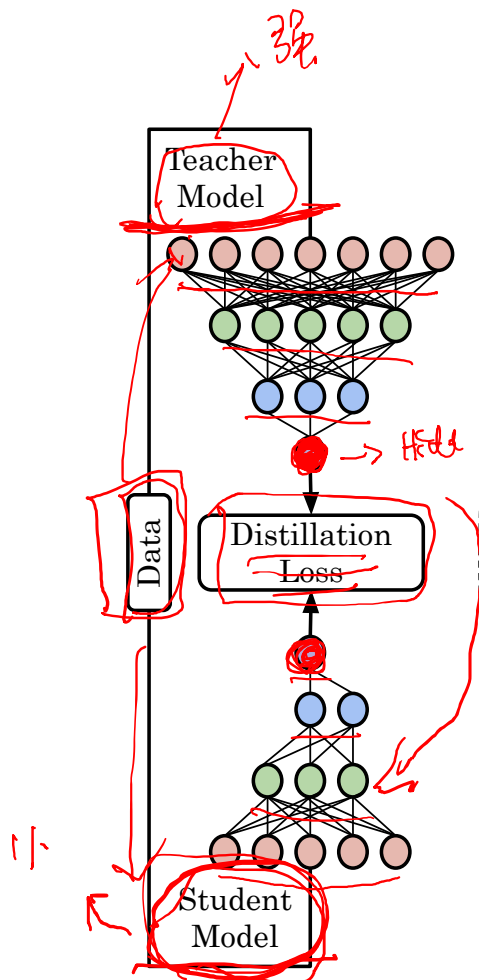


(e) block local

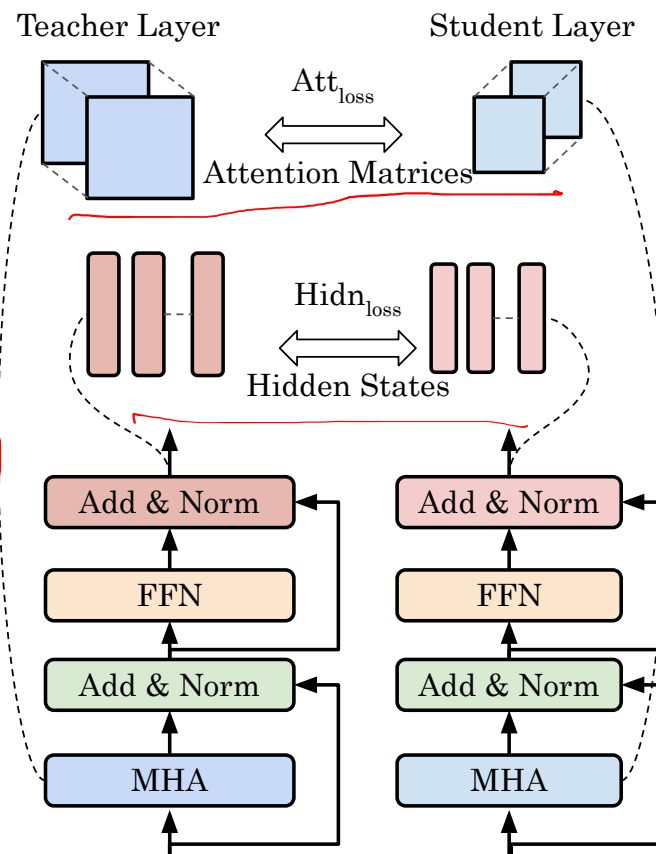


推理与压缩

- 剪枝
- 蒸馏
- 量化



(a) KD Overview



(b) Transformer Distillation



COT + Multi-agent → idea 1

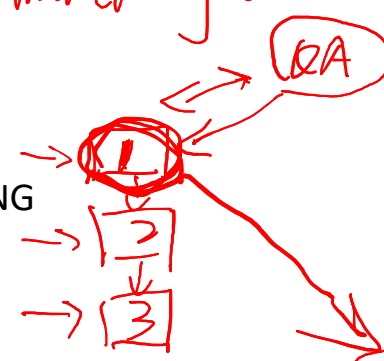
Socratic Method

□ → - 次推理

不同模型
不同角色
role-play

QA: prompt COT → LLaMA → output
A
COT... A

↓ multi-agent



<https://princeton-nlp.github.io/SocraticAI/>

TOT

- CRITIC: LARGE LANGUAGE MODELS CAN SELF-CORRECT WITH TOOL-INTERACTIVE CRITIQUING
<https://arxiv.org/pdf/2305.11738.pdf>
- INSTRUCTSCORE: Explainable Text Generation Evaluation with Finegrained Feedback
<https://arxiv.org/abs/2305.14282>
- Pinpoint, Not Criticize: Refining Large Language Models via Fine-Grained Actionable Feedback

推理: 认为 A 真, 所以 A 真

- Prompting Large Language Models With the Socratic Method

<https://arxiv.org/pdf/2303.08769.pdf>

A ✓

反思

前世记忆

→ boy → 知识

- SOCREVAL: Large Language Models with the Socratic Method for Reference-Free Reasoning Evaluation

<https://arxiv.org/pdf/2310.00074.pdf>

地图

演绎 新知

- SocraSynth: Socratic Synthesis for Reasoning and Decision Making

https://www.researchgate.net/publication/373753725_SocraSynth_Socratic_Synthesis_for_Reasoning_and_Decision_Making

IT → 300 G → COT → 反思

指导

Walton
Argument scheme

记忆

Critical question

