

第17章 场景题与开放题

场景题考察问题分析能力和工程经验，没有标准答案，重在思路

17.1 如何降低模型幻觉？

17.1.1 问题分析

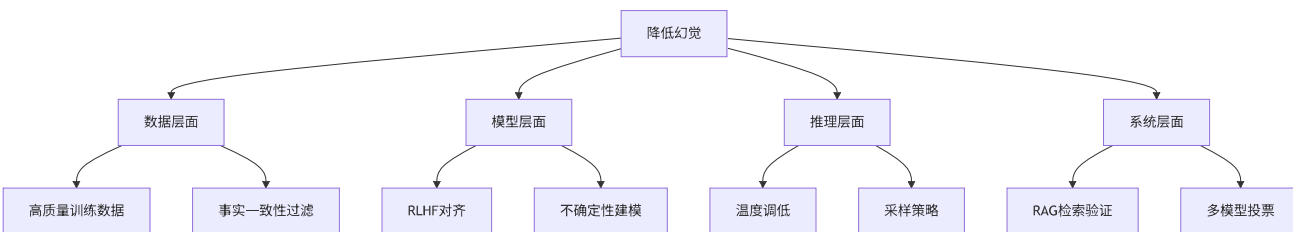
什么是幻觉？

- 模型编造不存在的事实
- 给出自信但错误的答案
- 前后矛盾的输出

幻觉的根本原因：

- 训练数据**：包含错误信息或偏见
- 模型机制**：语言建模目标不等于事实准确性
- 过度泛化**：从有限数据推广过度

17.1.2 解决方案矩阵



详细方案：

1. 训练阶段

```
def reduce_hallucination_in_training():  
    """  
    训练阶段的策略  
    """  
    strategies = {  
        "数据质量": {  
            "清洗": "移除低质量、错误数据",  
            "验证": "事实核查（维基百科、可靠来源）",  
            "多样性": "多来源验证，避免单一偏见"  
        },  
        "训练目标": {  
            "RLHF": "奖励模型惩罚幻觉行为",  
            "DPO": "偏好数据中标注幻觉为负例",  
            "Constitutional AI": "自我批评和改进"  
        },  
    }
```

```

        "不确定性建模": {
            "校准": "让模型知道自己不知道",
            "拒绝回答": "训练模型说 '我不知道'",
            "置信度": "输出答案的置信度分数"
        }
    }

    return strategies

```

2. 推理阶段

```

def reduce_hallucination_in_inference(query, model, knowledge_base):
    """
    推理阶段的策略
    """

    # 策略1: RAG (检索增强生成)
    relevant_docs = retrieve_from_knowledge_base(query, knowledge_base)

    prompt = f"""
    基于以下可靠来源回答问题。如果来源中没有相关信息，明确说"根据提供的信息无法回答"。

    来源: {relevant_docs}

    问题: {query}

    要求:
    1. 只使用来源中的信息
    2. 不要编造细节
    3. 不确定时明确指出
    """

    # 策略2: 多次采样 + 一致性检查
    n_samples = 5
    responses = []
    for _ in range(n_samples):
        response = model.generate(prompt, temperature=0.7)
        responses.append(response)

    # 检查一致性
    consistency_score = check_consistency(responses)

    if consistency_score < 0.8:
        # 一致性低，可能有幻觉
        return "我不太确定答案，建议查阅权威资料。"

    # 策略3: 自我验证
    final_response = most_common_response(responses)

```

```

verification_prompt = f"""
请验证以下答案是否符合常识和逻辑：

问题: {query}
答案: {final_response}

验证：
1. 是否合理？
2. 是否有明显错误？
3. 是否需要修正？
"""

verification = model.generate(verification_prompt, temperature=0.2)

return final_response, verification


def check_consistency(responses):
    """
    检查多个回答的一致性
    """
    from sentence_transformers import SentenceTransformer

    model = SentenceTransformer('all-MiniLM-L6-v2')
    embeddings = model.encode(responses)

    # 计算embedding的平均余弦相似度
    from sklearn.metrics.pairwise import cosine_similarity
    similarities = cosine_similarity(embeddings)

    # 上三角（不包括对角线）的平均相似度
    consistency = similarities[np.triu_indices_from(similarities, k=1)].mean()

    return consistency

```

3. 系统层面

```

class HallucinationGuard:
    """
    幻觉防护系统
    """
    def __init__(self, model, knowledge_base, fact_checker):
        self.model = model
        self.kb = knowledge_base
        self.fact_checker = fact_checker

    def generate_with_guard(self, query):

```

```

"""
带防护的生成
"""

# 1. 检测是否需要事实信息
if self.requires_factual_info(query):
    # 使用RAG
    response = self.generate_with_rag(query)
else:
    # 创造性任务，直接生成
    response = self.model.generate(query)

# 2. 事实核查
if self.contains_factual_claims(response):
    facts = self.extract_facts(response)
    verified_facts = self.fact_checker.verify(facts)

    # 标注未验证的事实
    response = self.annotate_unverified(response, verified_facts)

# 3. 不确定性标注
confidence = self.estimate_confidence(response)

if confidence < 0.7:
    response += "\n\n[注意：此答案的置信度较低，建议核实。]"

return response

def requires_factual_info(self, query):
    """
    判断查询是否需要事实信息
    """
    factual_keywords = [
        "什么时候", "多少", "哪里", "谁",
        "事实", "数据", "统计", "历史"
    ]
    return any(kw in query for kw in factual_keywords)

def extract_facts(self, text):
    """
    提取文本中的事实性陈述
    """
    # 使用NER、关系抽取等技术
    facts = []
    # ... 实现细节
    return facts

```

4. 评估和监控

```
def evaluate_hallucination_rate(model, test_dataset):
    """
    评估幻觉率
    """
    results = {
        "total": len(test_dataset),
        "hallucinations": 0,
        "unknown_rate": 0,
        "correct_rate": 0
    }

    for item in test_dataset:
        question = item["question"]
        ground_truth = item["answer"]

        response = model.generate(question)

        # 检查是否幻觉
        if "不知道" in response or "无法回答" in response:
            results["unknown_rate"] += 1
        elif verify_factual_correctness(response, ground_truth):
            results["correct_rate"] += 1
        else:
            results["hallucinations"] += 1

    # 计算比例
    for key in ["hallucinations", "unknown_rate", "correct_rate"]:
        results[key] = results[key] / results["total"]

    return results
```

17.1.3 最佳实践总结

阶段	方法	优先级	成本
训练	RLHF对齐	高	高
训练	高质量数据	高	高
推理	RAG	最高	中
推理	自我一致性	中	中
系统	事实核查	高	中
系统	多模型投票	低	高

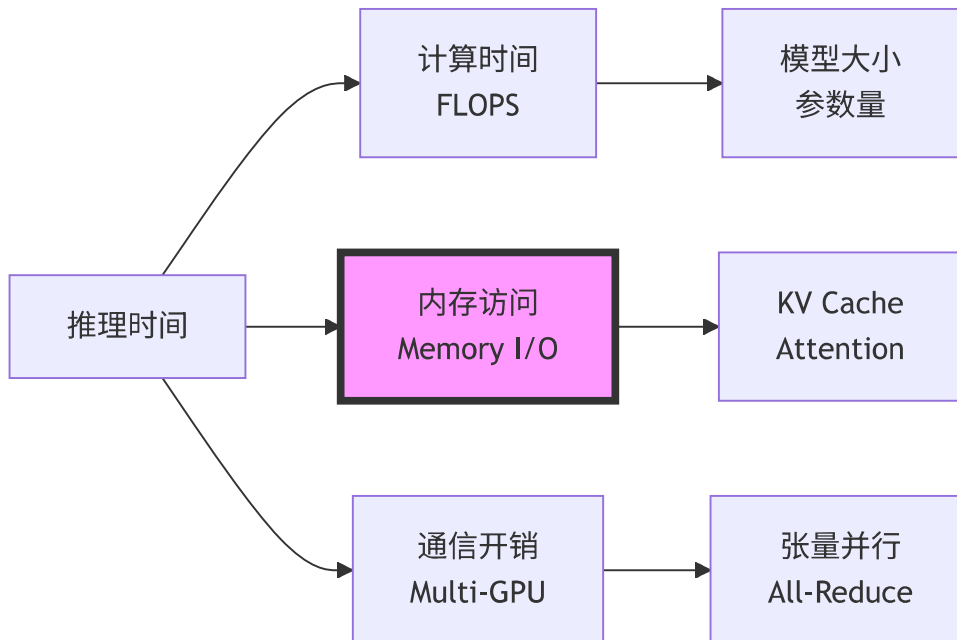
回答要点：

1. 承认幻觉是当前大模型的固有问题
2. 提出多层次的解决方案

3. 强调RAG作为最实用的方案
4. 提及评估和监控的重要性

17.2 如何提升模型推理速度？

17.2.1 瓶颈分析



关键观察：

- 对于大模型，瓶颈往往是**内存带宽**，不是计算
- 生成任务是自回归的，无法完全并行

17.2.2 优化策略

1. 模型层面

```
optimization_strategies = {  
    "量化": {  
        "INT8量化": {  
            "加速": "2-3x",  
            "精度损失": "<1%",  
            "实现": "使用GPTQ/AWQ",  
            "适用": "推理"  
        },  
        "INT4量化": {  
            "加速": "3-4x",  
            "精度损失": "1-3%",  
            "实现": "QLoRA",  
            "适用": "推理"  
        }  
    },  
    "剪枝": {
```

```

    "结构化剪枝": {
        "加速": "1.5-2x",
        "方法": "剪掉注意力头、FFN神经元",
        "适用": "推理"
    },
    "非结构化剪枝": {
        "加速": "有限（需要稀疏矩阵库）",
        "稀疏度": "50-90%",
        "适用": "特定硬件"
    }
},

"知识蒸馏": {
    "描述": "用小模型学习大模型",
    "加速": "5-10x (70B → 7B)",
    "精度损失": "5-15%",
    "成本": "需要训练"
}
}

```

2. 推理技巧

```

def optimize_inference():
    """
    推理层面的优化
    """
    optimizations = {
        "KV Cache": {
            "原理": "缓存已计算的Key和Value",
            "加速": "2-3x（长序列）",
            "代价": "增加显存占用",
            "必须使用": True
        },

        "Batch推理": {
            "原理": "合并多个请求一起处理",
            "加速": "N倍（batch_size=N时）",
            "挑战": "需要padding，处理不同长度",
            "吞吐量提升": "5-10x"
        },

        "连续Batching": {
            "原理": "动态组batch, iteration-level",
            "代表": "vLLM的continuous batching",
            "吞吐量提升": "10-20x",
            "推荐": "生产环境首选"
        },
    },

```

```

        "Speculative Decoding": {
            "原理": "小模型生成，大模型验证",
            "加速": "2-3x",
            "无损": "输出分布完全相同",
            "需要": "小模型+大模型"
        }
    }

    return optimizations

```

具体实现：Speculative Decoding

```

def speculative_decoding(large_model, small_model, input_ids,
                        k=5, max_length=100):
    """
    推测解码：小模型快速生成，大模型验证

    Args:
        large_model: 目标大模型
        small_model: 小模型
        k: 每次小模型生成k个token
    """
    current_ids = input_ids

    while len(current_ids[0]) < max_length:
        # 1. 小模型快速生成k个token
        draft_ids = small_model.generate(
            current_ids,
            max_new_tokens=k,
            do_sample=False # greedy
        )

        # 2. 大模型验证
        with torch.no_grad():
            large_logits = large_model(draft_ids).logits[:, -k-1:, :]
            large_probs = F.softmax(large_logits, dim=-1)

            small_logits = small_model(draft_ids).logits[:, -k-1:-1, :]
            small_probs = F.softmax(small_logits, dim=-1)

        # 3. 逐个验证token
        accepted_tokens = []
        for i in range(k):
            draft_token = draft_ids[0, -k+i]

            # 检查大模型是否接受这个token
            acceptance_prob = large_probs[0, i, draft_token] / small_probs[0, i, draft_t

```



```

        if random.random() < min(1, acceptance_prob):
            # 接受
            accepted_tokens.append(draft_token)
        else:
            # 拒绝, 从大模型采样新token
            new_token = torch.multinomial(large_probs[0, i], 1)
            accepted_tokens.append(new_token.item())
            break # 后续token都作废

# 4. 更新序列
current_ids = torch.cat([
    current_ids,
    torch.tensor([accepted_tokens], device=current_ids.device)
], dim=1)

# 5. 如果遇到EOS, 结束
if accepted_tokens[-1] == eos_token_id:
    break

return current_ids

```

3. 系统层面

```

class OptimizedInferenceSystem:
    """
    优化的推理系统
    """
    def __init__(self, model_path):
        # 1. 模型加载优化
        self.model = self.load_model_optimized(model_path)

        # 2. 编译优化
        self.model = torch.compile(self.model) # PyTorch 2.0+

        # 3. 使用高效推理引擎
        self.use_vllm = True

    def load_model_optimized(self, model_path):
        """
        优化的模型加载
        """
        from transformers import AutoModelForCausalLM

        model = AutoModelForCausalLM.from_pretrained(
            model_path,
            torch_dtype=torch.float16, # FP16
            device_map="auto",         # 自动设备分配
            low_cpu_mem_usage=True,    # 降低CPU内存

```

```

    )

    return model

def generate_batch(self, prompts, max_length=100):
    """
    批量生成
    """
    if self.use_vllm:
        # 使用vLLM (推荐)
        outputs = self.vllm_generate(prompts, max_length)
    else:
        # HuggingFace标准方式
        inputs = self.tokenizer(prompts, return_tensors="pt", padding=True)
        outputs = self.model.generate(**inputs, max_length=max_length)

    return outputs

def benchmark_inference_speed(model, input_length=100, output_length=50, batch_size=1):
    """
    基准测试推理速度
    """
    import time

    # 准备输入
    input_ids = torch.randint(0, 50000, (batch_size, input_length))

    # Warmup
    for _ in range(5):
        _ = model.generate(input_ids, max_new_tokens=10)

    # 测试
    start = time.time()
    num_runs = 10

    for _ in range(num_runs):
        outputs = model.generate(input_ids, max_new_tokens=output_length)

    end = time.time()

    # 计算指标
    total_time = end - start
    avg_latency = total_time / num_runs # 平均延迟 (秒)
    tokens_per_second = (batch_size * output_length * num_runs) / total_time

    print(f"平均延迟: {avg_latency:.2f}秒")
    print(f"吞吐量: {tokens_per_second:.1f} tokens/s")
    print(f"Batch大小: {batch_size}")

```

```
return {
    "latency": avg_latency,
    "throughput": tokens_per_second
}
```

4. 硬件层面

```
hardware_considerations = {
    "GPU选择": {
        "推理": "A100 (80GB) > A10 > T4",
        "考虑": "显存容量、带宽、价格",
        "量化后": "更小的GPU也可以（如4090）"
    },

    "Flash Attention": {
        "加速": "2-4x（长序列）",
        "要求": "A100/H100（计算能力8.0+）",
        "降低": "显存占用50%+"
    },

    "张量并行": {
        "适用": "单卡放不下的大模型",
        "通信": "NVLink > PCIe",
        "效率": "2卡：70-80%， 4卡： 60-70%"
    }
}
```

17.2.3 优化效果对比

方法	加速比	精度损失	实现难度	推荐指数
KV Cache	2-3x	无	简单	★★★★★
INT8量化	2x	<1%	简单	★★★★★
vLLM	10-20x	无	简单	★★★★★
Flash Attention	2-4x	无	中等	★★★★★
Speculative Decoding	2-3x	无	复杂	★★★
知识蒸馏	5-10x	5-15%	很复杂	★★★

17.2.4 回答要点

面试官： "如何提升模型推理速度？ "

回答框架：

- 1. 【瓶颈分析】

- 先分析瓶颈：内存带宽 vs 计算
- 对于大模型，通常是内存瓶颈

2. 【分层方案】

- 模型层面：量化（最实用）、剪枝、蒸馏
- 推理技巧：KV Cache（必须）、Batch、Speculative Decoding
- 系统层面：vLLM、Flash Attention
- 硬件层面：选择合适的GPU

3. 【权衡】

- 量化：精度 vs 速度
- Batch：延迟 vs 吞吐量
- 蒸馏：准确性 vs 速度

4. 【实践经验】（如果有）

- 具体数字：从Xms降到Yms
- 使用的技术组合
- 遇到的坑和解决方案

5. 【推荐方案】

- 优先级：KV Cache > 量化 > vLLM > Flash Attention
- 生产环境：vLLM + INT8量化

17.3 如何处理超长文本？

17.3.1 问题定义

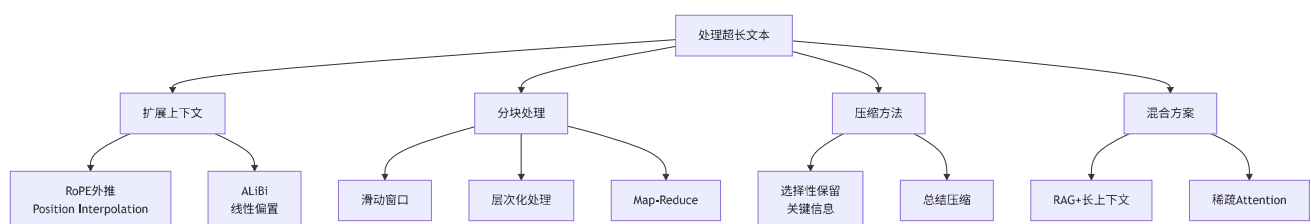
什么是超长文本？

- 超过模型上下文长度限制（如4K, 8K）
- 例子：长篇论文、完整书籍、法律文档

挑战：

1. **技术限制**：Attention复杂度 $O(n^2)$
2. **显存限制**：KV Cache随序列长度线性增长
3. **性能下降**：“迷失在中间”（Lost in the Middle）

17.3.2 解决方案



方案1: 分块+滑动窗口

```

class SlidingWindowProcessor:
    """
    滑动窗口处理超长文本
    """

    def __init__(self, model, max_length=2048, stride=1024):
        self.model = model
        self.max_length = max_length
        self.stride = stride # 步长, 窗口重叠

    def process_long_document(self, document, task="summarize"):
        """
        处理超长文档

        Args:
            document: 长文档
            task: 任务类型 (summarize, qa, etc.)
        """
        chunks = self.split_into_windows(document)

        if task == "summarize":
            return self.summarize_windows(chunks)
        elif task == "qa":
            return self.qa_over_windows(chunks)

    def split_into_windows(self, text):
        """
        分割成重叠的窗口
        """
        tokens = self.tokenizer.encode(text)
        windows = []

        start = 0
        while start < len(tokens):
            end = min(start + self.max_length, len(tokens))
            window = tokens[start:end]
            windows.append(self.tokenizer.decode(window))

            if end == len(tokens):
                break

            start += self.stride

        return windows

    def summarize_windows(self, windows):
        """
        总结窗口内容
        """

```

```

summaries = []

# 第一轮：总结每个窗口
for window in windows:
    prompt = f"请简要总结以下内容： \n{window}\n\n总结： "
    summary = self.model.generate(prompt)
    summaries.append(summary)

# 第二轮：合并总结
if len(summaries) > 1:
    combined = "\n\n".join(summaries)

    if len(combined) > self.max_length:
        # 递归处理
        return self.summarize_windows(
            self.split_into_windows(combined)
        )
    else:
        final_prompt = f"""
        以下是文档各部分的总结，请给出整体总结：

        {combined}

        整体总结：
        """
        return self.model.generate(final_prompt)

return summaries[0]

```

方案2: Map-Reduce模式

```

def map_reduce_long_document(document, model, task="summarize"):
    """
    Map-Reduce处理长文档

    Map：分块处理
    Reduce：合并结果
    """
    # Map阶段：分块
    chunk_size = 2000 # 字符数
    chunks = [document[i:i+chunk_size]
                for i in range(0, len(document), chunk_size)]

    # Map阶段：并行处理每个块
    from concurrent.futures import ThreadPoolExecutor

    def process_chunk(chunk):
        prompt = f"总结： {chunk}\n\n摘要： "

```

```

        return model.generate(prompt)

with ThreadPoolExecutor(max_workers=4) as executor:
    chunk_summaries = list(executor.map(process_chunk, chunks))

# Reduce阶段: 合并
combined_summary = "\n".join(chunk_summaries)

final_prompt = f"""
以下是文档分段的摘要，请整合成一个连贯的总结：

{combined_summary}

整体总结：
"""

final_summary = model.generate(final_prompt)

return final_summary

```

方案3: 稀疏Attention

```

class SparseAttention:
    """
    稀疏注意力: 只关注部分token
    """

    def __init__(self, pattern="strided"):
        self.pattern = pattern

    def create_sparse_mask(self, seq_len, block_size=64):
        """
        创建稀疏mask

        Args:
            seq_len: 序列长度
            block_size: 块大小

        Returns:
            mask: (seq_len, seq_len) 稀疏mask
        """
        mask = torch.zeros(seq_len, seq_len)

        if self.pattern == "strided":
            # Strided attention: 每隔几个位置关注一次
            stride = 128
            for i in range(seq_len):
                # 局部注意力 (前后block_size个token)
                mask[i, max(0, i-block_size):i+1] = 1

```

```

        # 跨步注意力
        for j in range(0, seq_len, stride):
            mask[i, j] = 1

    elif self.pattern == "fixed":
        # Fixed attention: 所有token都关注前几个token
        mask[:, :block_size] = 1

    # 加上局部注意力
    for i in range(seq_len):
        mask[i, max(0, i-block_size):i+1] = 1

    return mask

# Longformer风格的注意力
def longformer_attention(Q, K, V, window_size=256):
    """
    Longformer: 局部窗口 + 全局token

    时间复杂度:  $O(n * w)$  其中w是窗口大小
    """
    seq_len = Q.size(1)

    # 1. 局部窗口注意力
    local_mask = create_sliding_window_mask(seq_len, window_size)
    local_attn = masked_attention(Q, K, V, local_mask)

    # 2. 全局注意力 (特殊token, 如[CLS])
    global_token_ids = [0] # 第一个token

    for idx in global_token_ids:
        # 全局token关注所有token
        local_mask[idx, :] = 1
        # 所有token关注全局token
        local_mask[:, idx] = 1

    output = masked_attention(Q, K, V, local_mask)

    return output

```

方案4: 层次化处理

```

class HierarchicalProcessor:
    """
    层次化处理: 先处理段落, 再处理文档
    """

```



```

def __init__(self, model):
    self.model = model

def process(self, document):
    """
    层次化处理文档
    """
    # 1. 分割成段落
    paragraphs = self.split_paragraphs(document)

    # 2. 每个段落生成embedding或摘要
    paragraph_repr = []
    for para in paragraphs:
        summary = self.model.generate(f"一句话总结: {para}")
        paragraph_repr.append(summary)

    # 3. 构建段落级别的上下文
    document_structure = "\n".join([
        f"段落{i+1}: {summary}"
        for i, summary in enumerate(paragraph_repr)
    ])

    # 4. 回答问题时, 先找相关段落
    def answer_question(question):
        # 找最相关的段落
        relevant_paragraphs = self.retrieve_paragraphs(
            question, paragraphs, paragraph_repr
        )

        # 在相关段落中找答案
        context = "\n\n".join(relevant_paragraphs)

        answer = self.model.generate(f"""
        上下文: {context}

        问题: {question}

        答案:
        """)

        return answer

    return answer_question

```

方案5: 混合方案 (推荐)

```

class HybridLongContextSolution:
    """

```

混合方案：RAG + 长上下文模型

```
"""
def __init__(self, model, vectorstore, long_context_limit=128000):
    self.model = model # 假设支持128K上下文
    self.vectorstore = vectorstore
    self.long_context_limit = long_context_limit

def process(self, document, query):
    """
    处理超长文档的查询
    """
    doc_length = len(self.tokenizer.encode(document))

    if doc_length <= self.long_context_limit:
        # 能放进上下文，直接使用
        return self.direct_inference(document, query)
    else:
        # 太长，使用RAG
        return self.rag_inference(document, query)

def direct_inference(self, document, query):
    """
    直接推理（文档放入上下文）
    """
    prompt = f"""
    文档：
    {document}

    问题： {query}

    答案：
    """
    return self.model.generate(prompt)

def rag_inference(self, document, query):
    """
    RAG推理（检索相关片段）
    """
    # 1. 分块并索引
    chunks = self.chunk_document(document)
    self.vectorstore.add_documents(chunks)

    # 2. 检索相关片段
    relevant_chunks = self.vectorstore.similarity_search(query, k=10)

    # 3. 用长上下文模型处理检索到的片段
    context = "\n\n".join([chunk.page_content for chunk in relevant_chunks])

    prompt = f"""
```

相关片段:

{context}

问题: {query}

答案:

"""

return self.model.generate(prompt)

17.3.3 方案对比

方案	适用场景	优点	缺点	实现难度
滑动窗口	顺序阅读、总结	简单、保持连续性	可能丢失全局信息	★
Map-Reduce	并行处理、总结	可并行、可扩展	需要多次调用	★★
RAG	问答、检索	最实用、精准	依赖检索质量	★★
稀疏Attention	需要全文理解	理论上完美	需要定制模型	★★★★
长上下文模型	任意任务	最直接	成本高、可能迷失	★
混合方案	生产环境	效果最好	稍复杂	★★★

17.3.4 回答要点

面试官: "如何处理超长文本?"

回答框架:

1. 【明确场景】

- 超长是多长? (10K? 100K? 1M?)
- 任务类型? (问答、总结、分类?)

2. 【分层方案】

基础方案:

- 分块+滑动窗口 (简单任务)
- Map-Reduce (总结任务)
- RAG (问答任务) ← 最实用

进阶方案:

- 稀疏Attention (需要定制)
- 层次化处理 (结构化文档)
- 长上下文模型 (GPT-4-128K)

3. 【权衡】

- 准确性 vs 效率
- 全局理解 vs 局部精确
- 实现复杂度 vs 效果

4. 【推荐】

- 10K以内：直接用长上下文模型
- 10-100K：RAG（最实用）
- 100K+：混合方案（RAG + 长上下文）

17.4 本章小结

本章介绍了大模型面试中的经典场景题：

✅ **降低幻觉**：RLHF + RAG + 自我一致性 ✅ **提升速度**：量化 + KV Cache + vLLM ✅ **超长文本**：RAG（最实用）+ 长上下文

回答场景题的框架：

1. **问题分析**：明确场景、约束、目标
2. **多维方案**：数据、模型、系统、硬件
3. **权衡取舍**：成本、效果、实现难度
4. **最佳实践**：推荐方案+优先级
5. **实际经验**（如果有）：具体项目经验

准备策略：

- 每个问题准备3-5分钟的回答
- 准备具体数字和案例
- 了解最新的技术进展
- 能画图解释复杂方案

下一章预告： 第18章将分析各大公司的面试特点和准备策略。