

# 大模型推理介绍



ZOMI

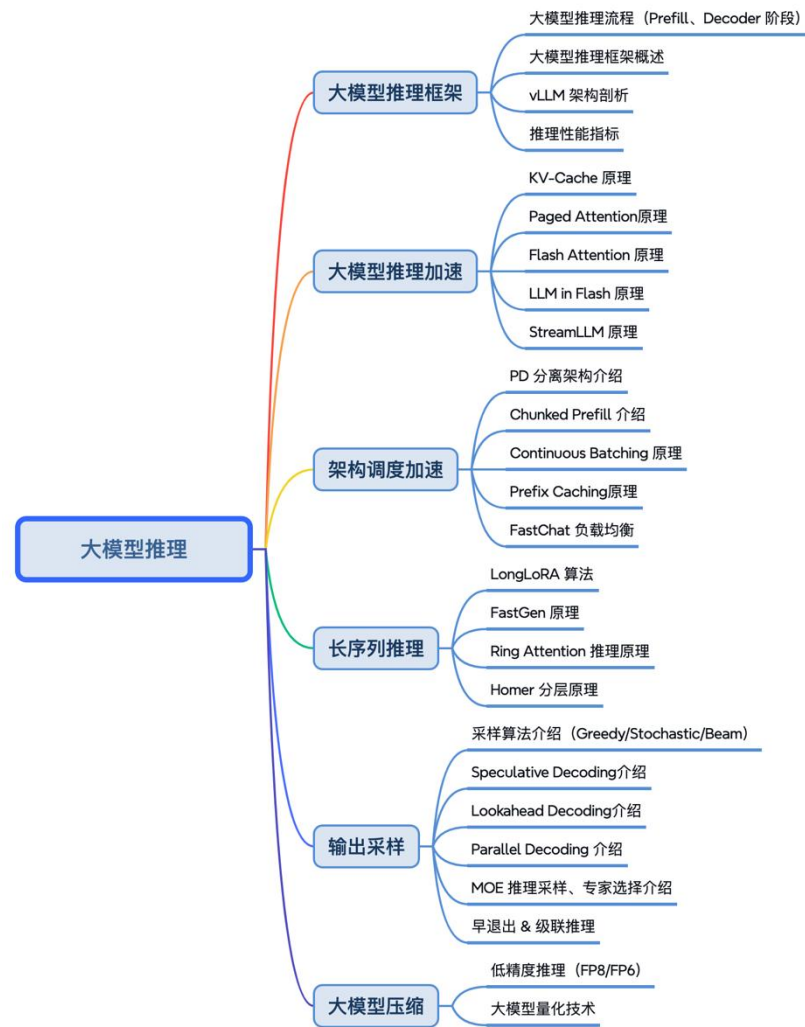
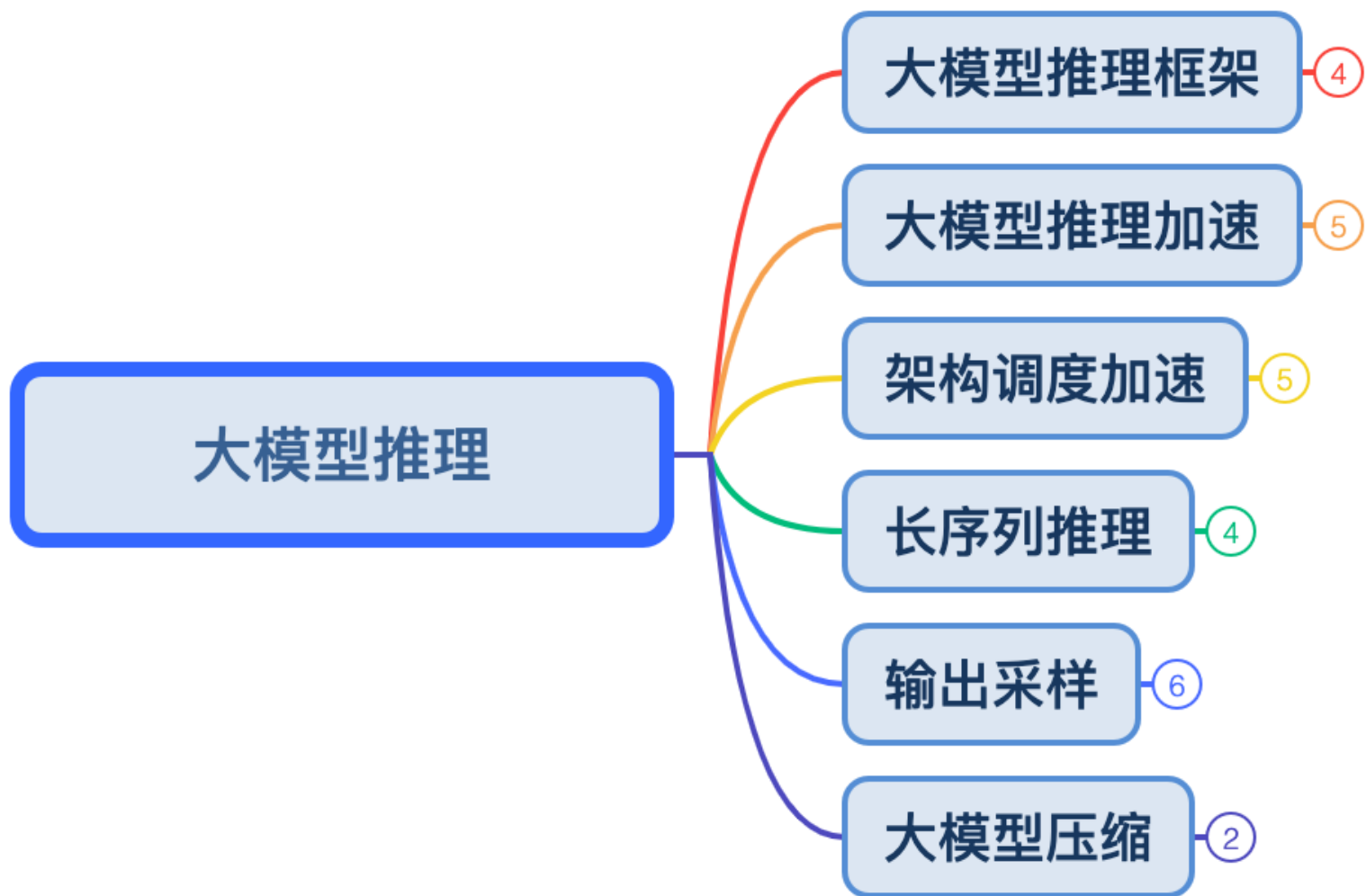
KV-cache  
transfer

# 思考

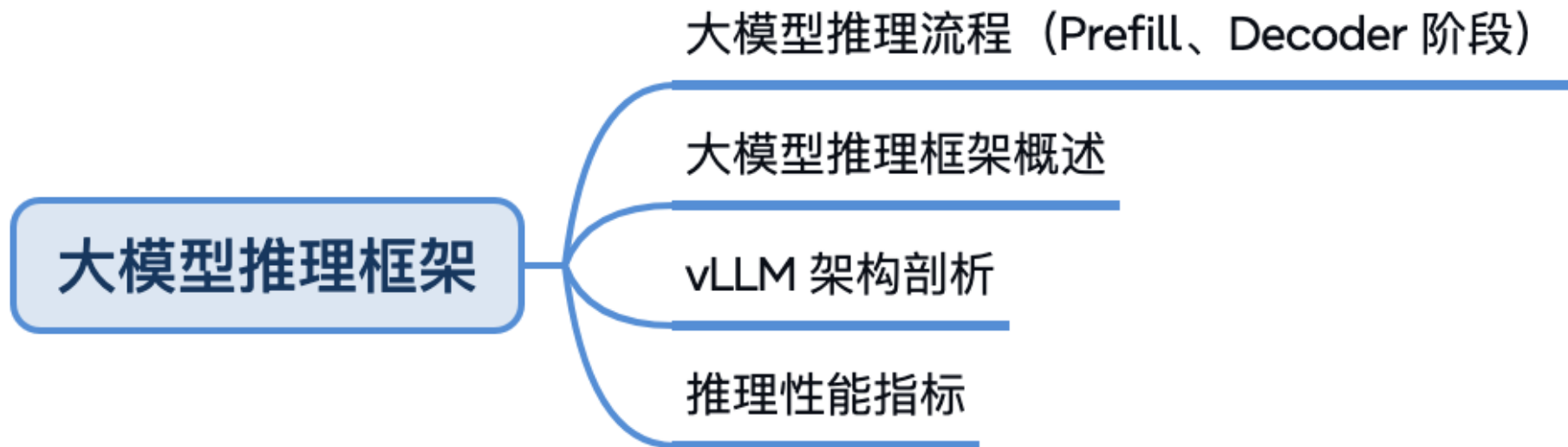
- 大模型推理框架跟传统推理框架，有什么不同？
  - 大模型推理框架只推理大模型，聚焦 Transformer 架构
  - 大模型推理要包含服务化，主要在云端使用
  - ??? 还有吗？那么架构上有什么区别？



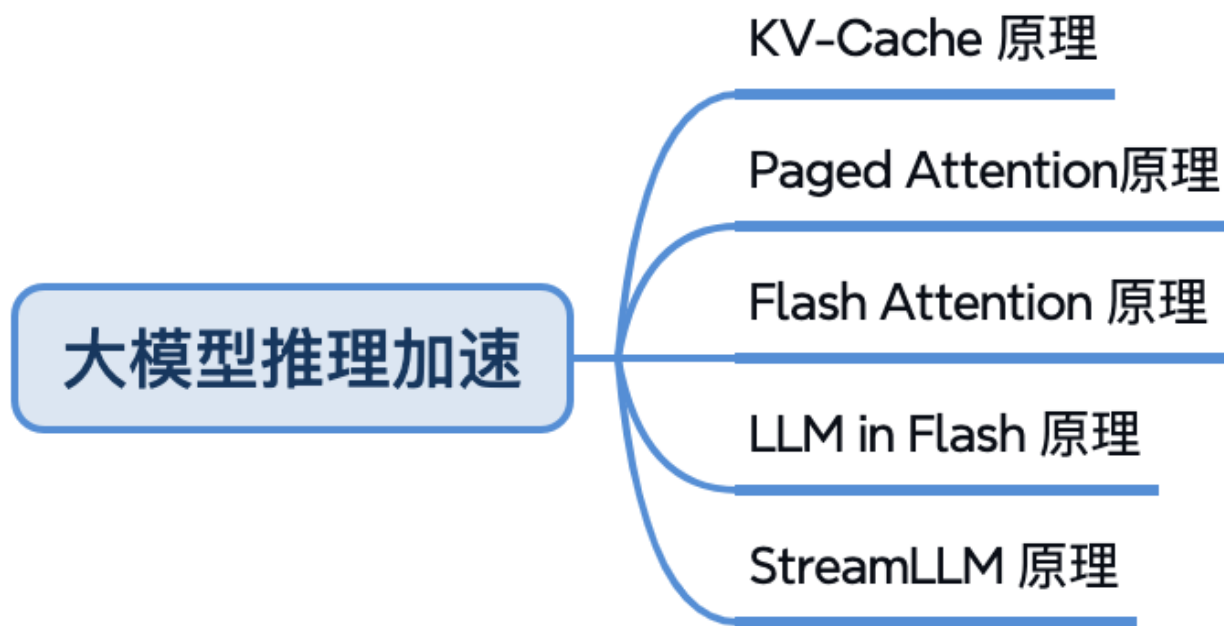
# 目录



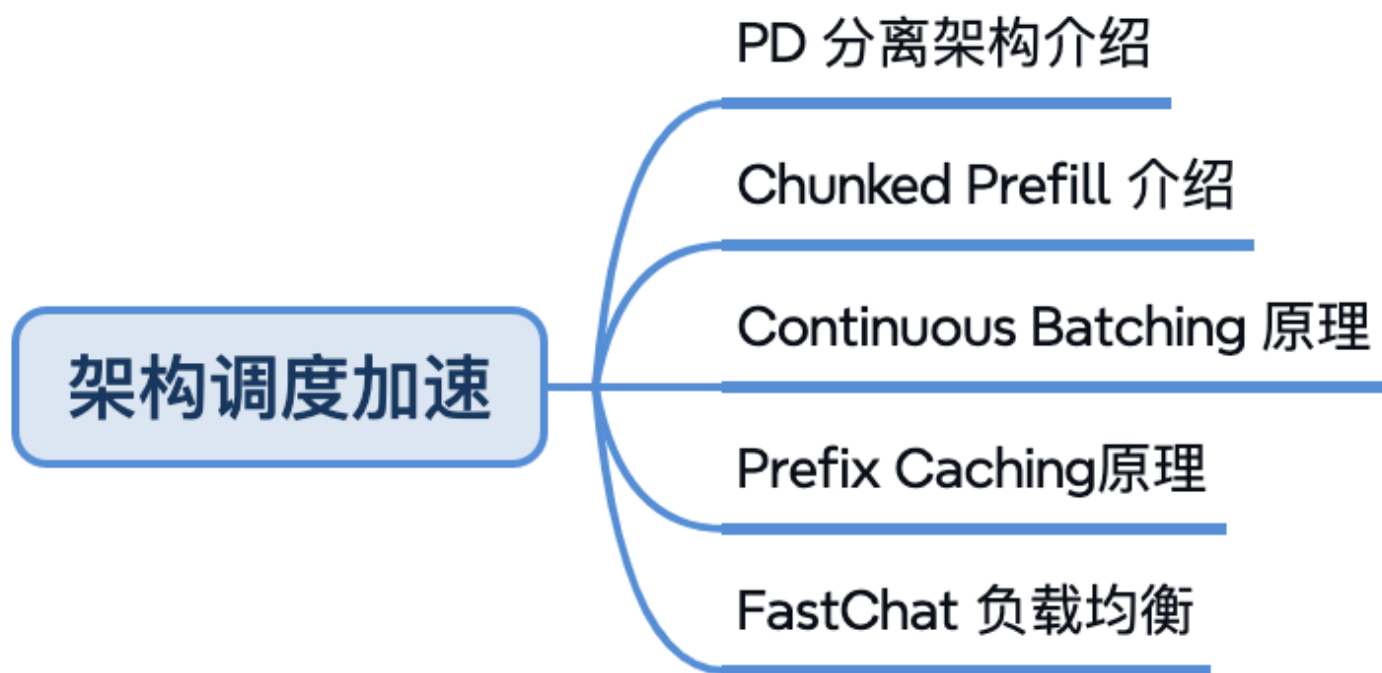
# 大模型推理



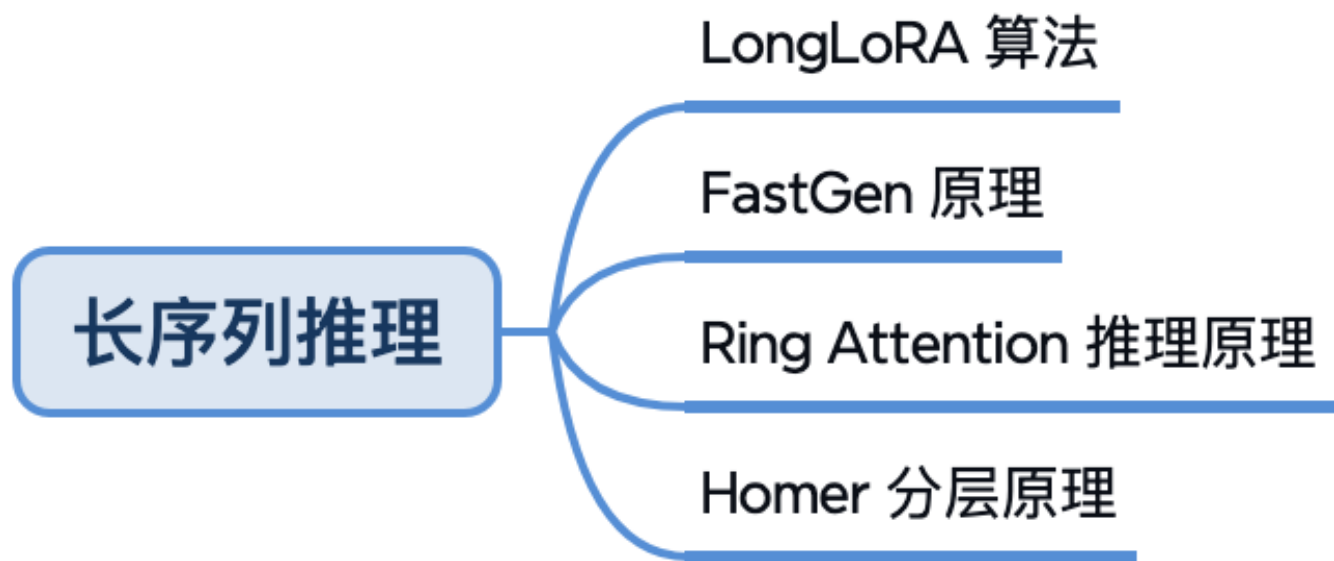
# 大模型推理



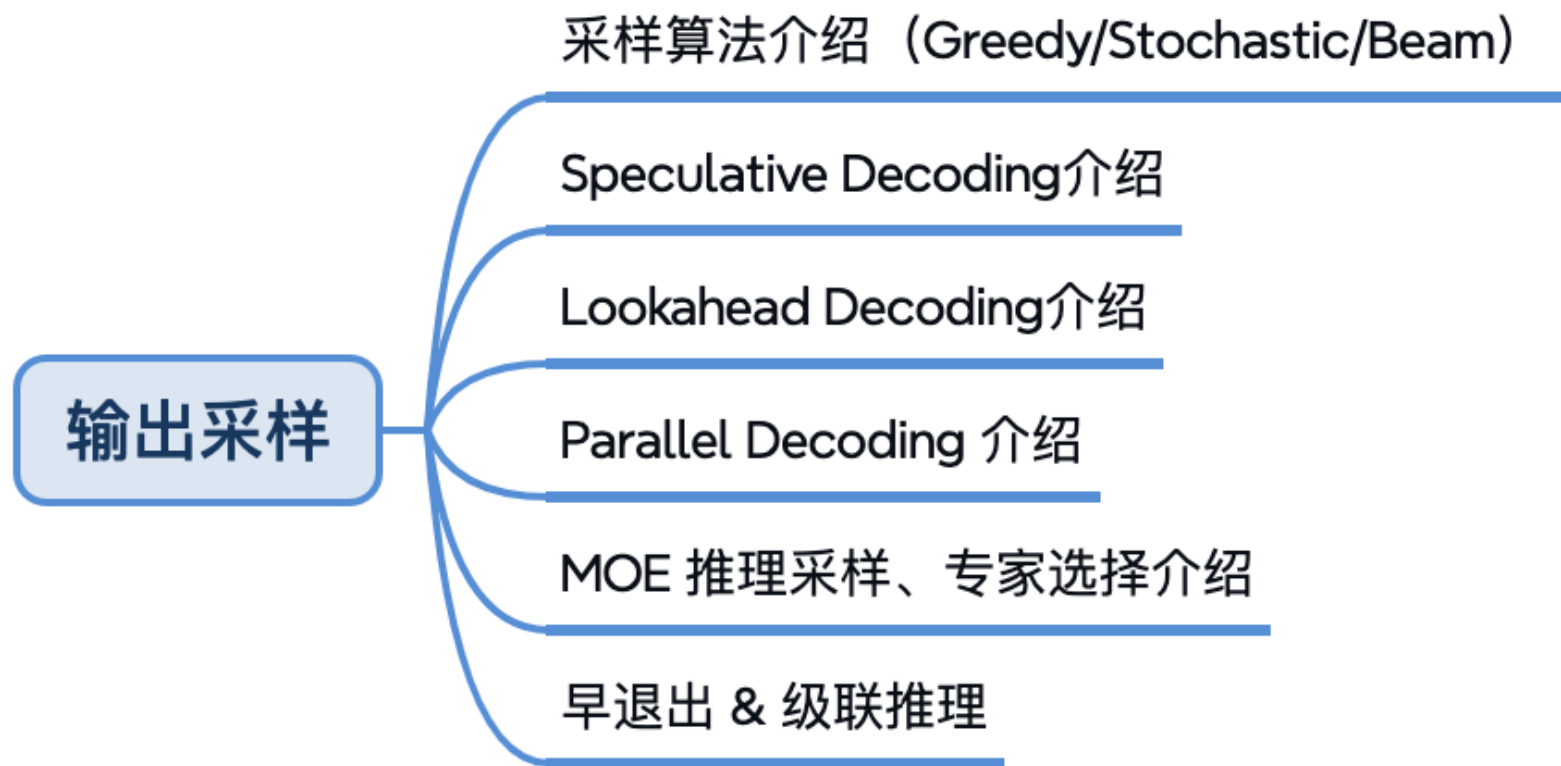
# 大模型推理



# 大模型推理



# 大模型推理





# 大模型推理





# Thank you

把AI系统带入每个开发者、每个家庭、  
每个组织，构建万物互联的智能世界

Bring AI System to every person, home and  
organization for a fully connected,  
intelligent world.

Copyright © 2024 XXX Technologies Co., Ltd.  
All Rights Reserved.

The information in this document may contain predictive statements including, without limitation, statements regarding the future financial and operating results, future product portfolio, new technology, etc. There are a number of factors that could cause actual results and developments to differ materially from those expressed or implied in the predictive statements. Therefore, such information is provided for reference purpose only and constitutes neither an offer nor an acceptance. XXX may change the information at any time without notice.



GitHub <https://github.com/chenzomi12/AIFoundation>