

RAG 技术详解与实践应用

第0讲：课程介绍

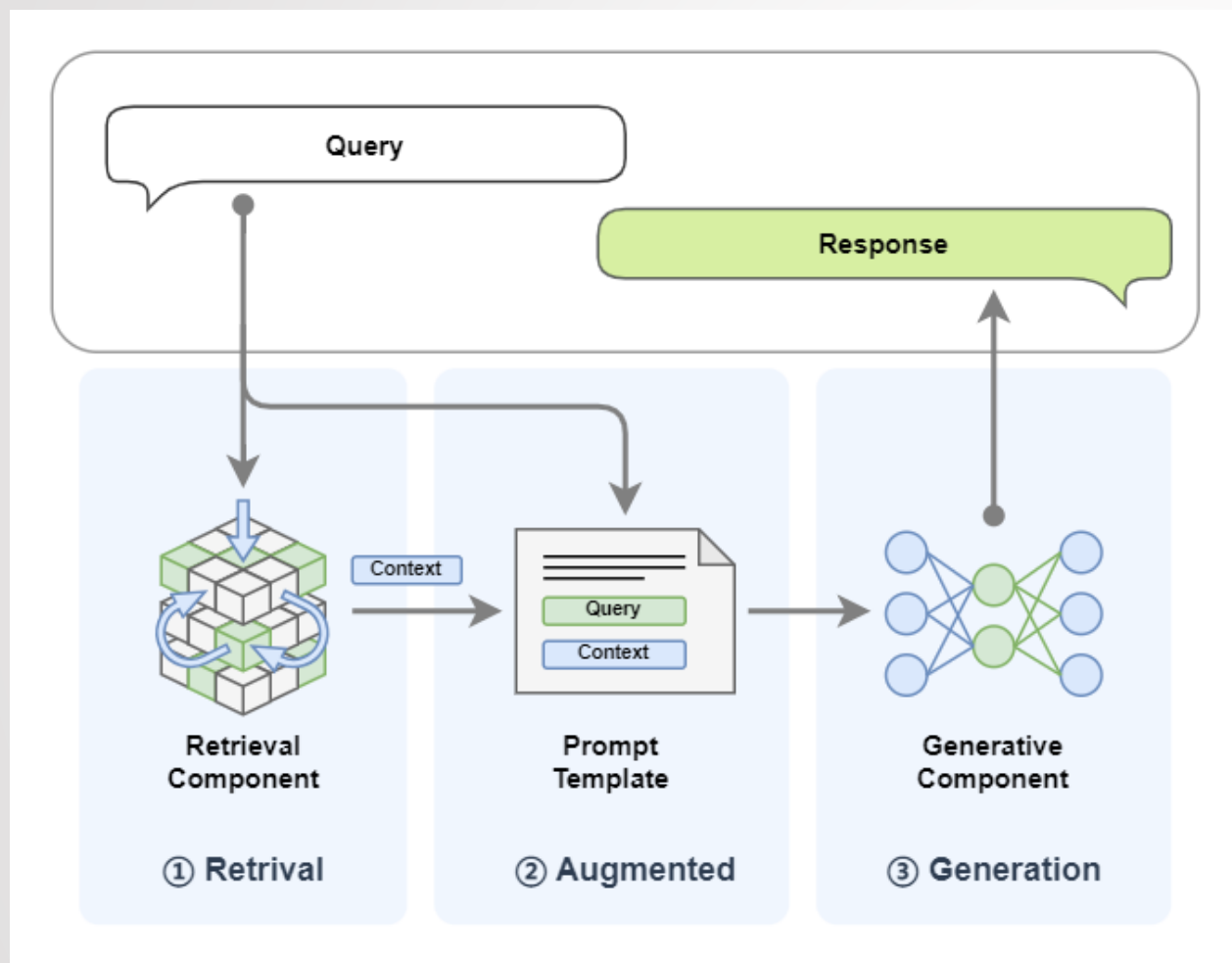


目录

-  1. 课程简介
-  2. 讲师介绍
-  3. 选题初衷
-  4. 课程导览



课程简介



以RAG(Retrieval-Augmented Generation)为核心，
从基础的检索与生成策略入手，
深入探索效果优化与性能提升，
逐步拓展多模态能力，
强大的智能Agent全面赋能，
构建可落地的企业级知识库。
通过大量的实战案例，
专业的导师手把手教学，
获得从原型开发到系统工程化的完整能力成长，
助力您斩获心仪的offer！



目录

-  1. 课程简介
-  2. 讲师介绍
-  3. 选题初衷
-  4. 课程导览



王志宏 商汤科技大装置事业群研发总监



- 商汤的私有化AI综合解决方案研发负责人，并主导开源项目 LazyLLM 的技术研发与生态建设
- 前商汤自研的深度学习训练框架SenseParrots的研发负责人
- 荣获开源中国源创会2024年度技术领航者称号
- 深耕 AI 领域多年，具备丰富的企业级RAG 实践经验，推动数十家企业实现 AI 应用落地，累计为企业创造了数千万元的商业价值



目录

-  1. 课程简介
-  2. 讲师介绍
-  3. 选题初衷
-  4. 课程导览



大模型时代的技术选型



技术方向	主要技能	学习难度	算力门槛	实用性
预训练 Pre-training	深度学习基础：反向传播、优化算法、梯度下降 自然语言处理：词向量、Transformer、Attention 机制 训练框架：PyTorch 大规模分布式训练：分布式计算、各种并行计算	★★★★	★★★★ ★★	★
微调 Fine-tuning	监督学习：分类、生成、问答任务建模和数据清洗 模型训练与优化：调参，过拟合处理 训练框架：PyTorch	★★	★★	★★
RAG	文本处理：文本读取，解析和切分策略 向量化检索：词向量化，向量数据库 模型使用：大模型，多模态模型，... 应用部署：服务化、API 接口开发，界面开发	★	★	★★★★ ★★
智能体 Agent	任务规划与调度：多任务管理、任务拆解 环境交互：API 集成、插件开发	★★	★	★★★
强化学习 RL	强化学习：学会一个强化学习的训练框架 算法定制：制定目标、Reward模型和Policy模型	★★★★ ★	★★★★	★★

RAG 的应用落地性最强，学习成本相对可控，因此是当前**最实用、最具商业价值**的方案。



	从 0 到 1 手搓 RAG	基于已有框架搭建并逐步优化效果
方案成熟度	对大部分人而言，稳定性和性能都不太好	稳定性和性能都比较好
开发周期	开发周期长，工作量大	快速搭建 MVP，进行效果验证
学习门槛	对代码能力要求较高	完成基础目标所需要的代码能力较低
定制化程度	灵活性高，可任意定制	不同框架的灵活度不同，但定制化程度均受限
学习重心	深入理解 RAG 工作原理	可快速聚焦优化检索效果和生成质量
适用用户群体	1. 想锻炼代码能力的人 2. 觉得已有框架不能满足需求的群体	1. RAG初学者 2. 企业开发者

结论：对于初学者而言，**基于已有框架**是明智的选择：

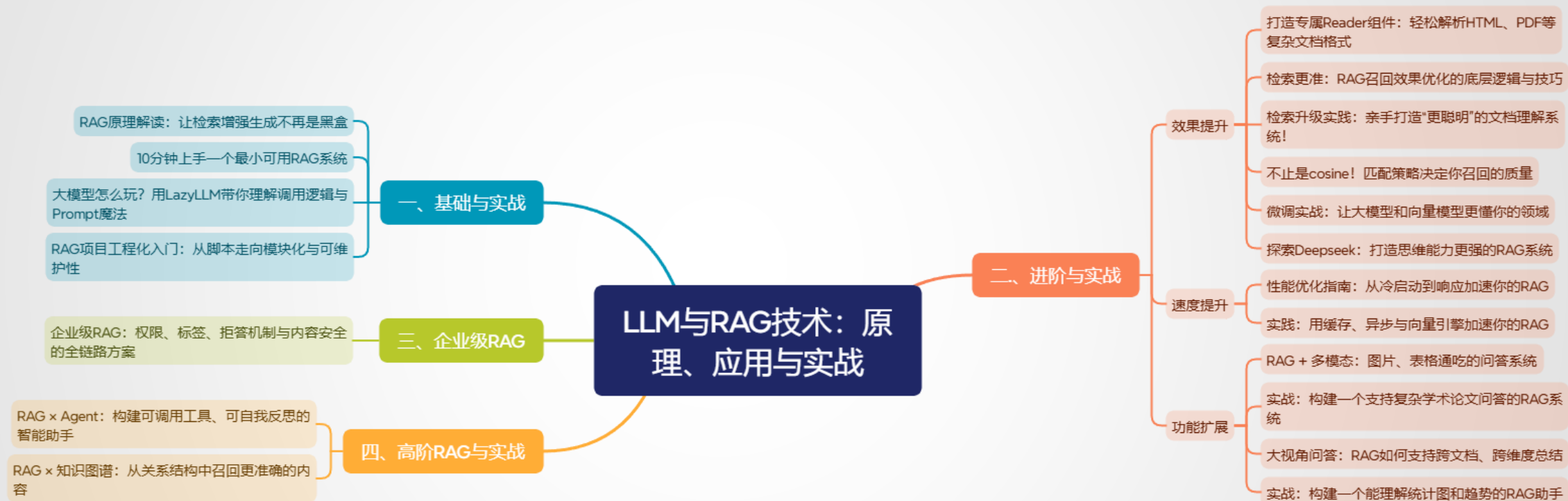
- 1.快速落地，先见成效：**已有框架提供了完备的检索、生成、调用接口，大大缩短开发周期。
- 2.成熟度高，性能优化好：**框架解决了很多复杂的工程问题，避免重复“造轮子”。
- 3.聚焦业务，注重成效：**专注应用场景的业务逻辑优化，而非底层细节，聚焦关键目标。



目录

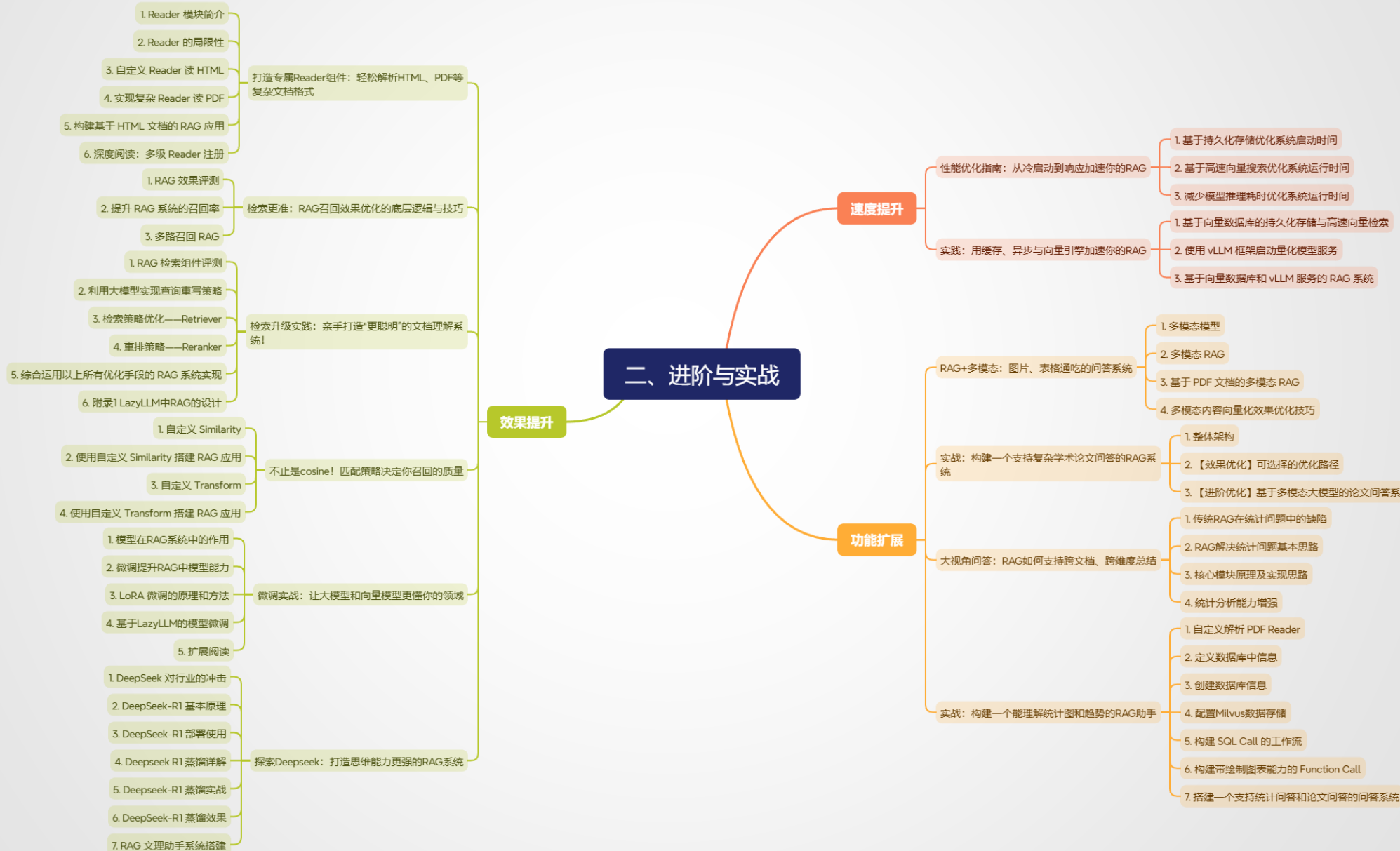
-  1. 课程简介
-  2. 讲师介绍
-  3. 选题初衷
-  4. 课程导览

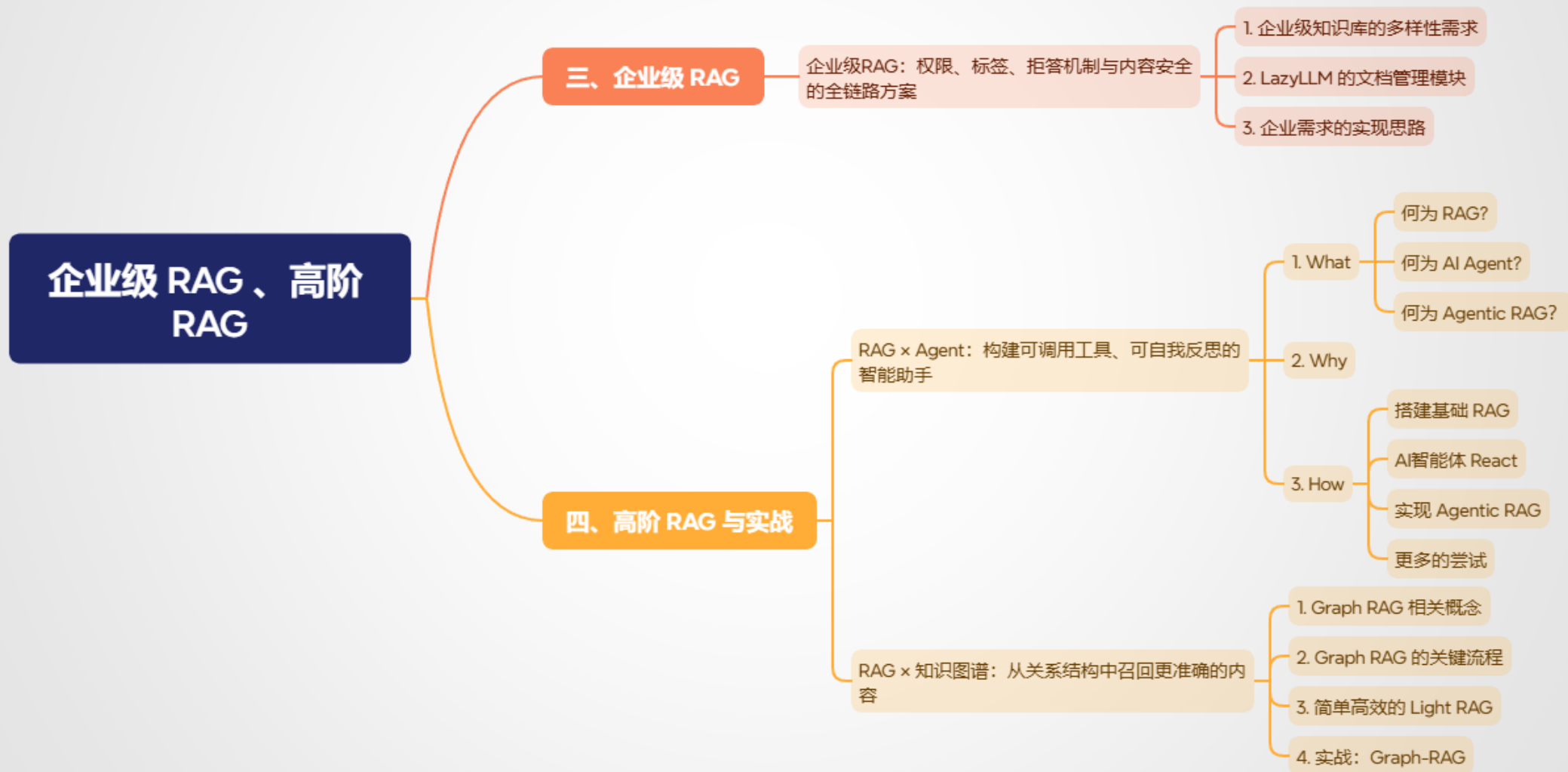






课程导览





Q&A

1. 为什么预训练的实用性只有一颗星
2. 课程上线节奏是什么样的



感谢聆听

Thanks for Listening