

QCon
全球软件开发大会

Serverless助力大语言模型工程化实践

演讲人：姬军翔

亚马逊科技 / 高级解决方案架构师

目录

01

LMSI模型

02

LLM应用案例分析

03

Serverless最佳实践

04

回顾和总结

多模态

关于 AI 的高频问题 都能在这里找到答案

RAG

AI 智驾

从大模型变革之路到高效“炼丹”指南

扫码领取你的智囊团

成本优化实践

AI Native 产品创新
与技术落地



咨询购票



查看详情

关于我



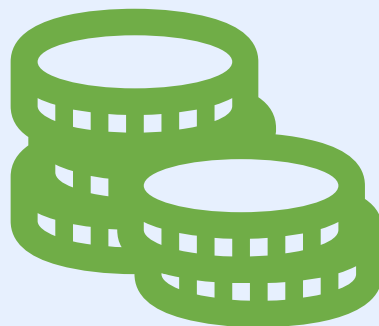
01 LMSI

语言模型系统接口模型

● 构建大语言模型应用的常见挑战



准确度，性能不达标



训练和部署昂贵



合规，构建复杂

一些应对方法

准确度，性能不达标



提示工程

部署成本高



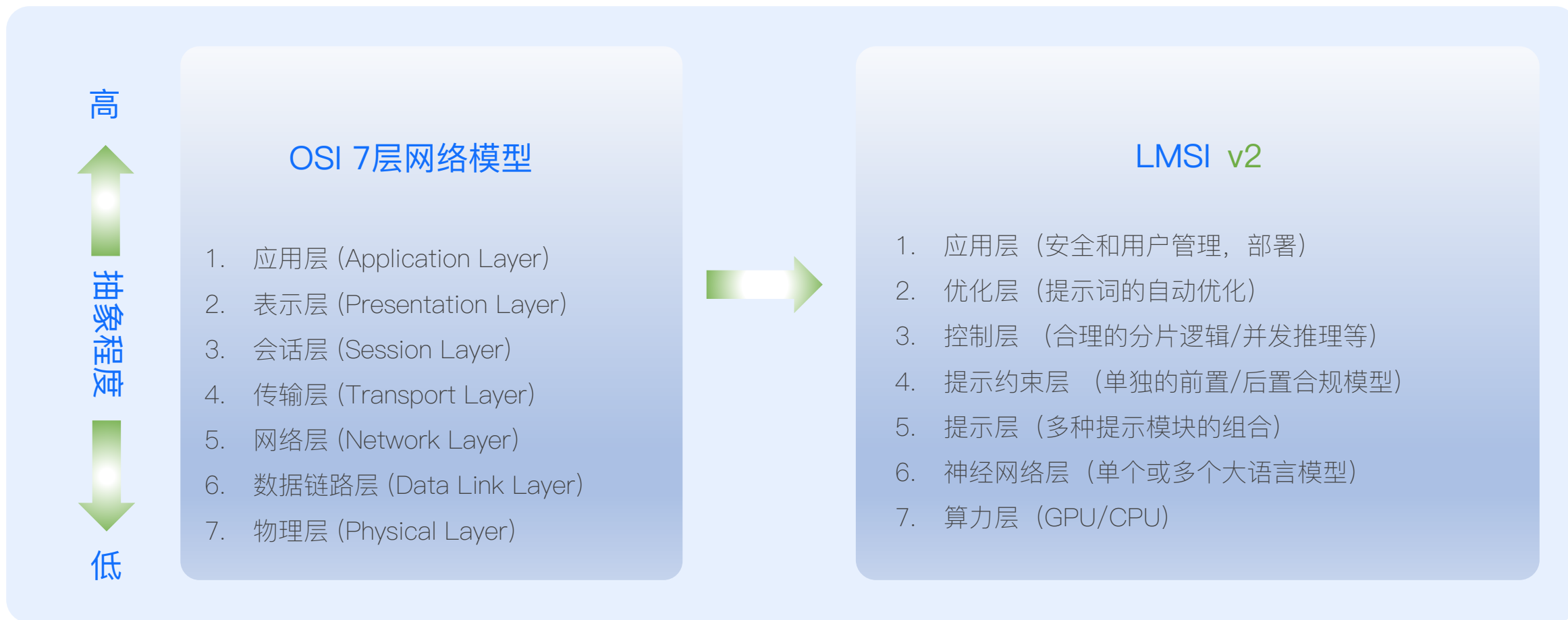
使用托管模型

构建复杂



Serverless/框架

语言模型系统接口模型(LMSI)



<https://www.twosigma.com/articles/a-guide-to-large-language-model-abstractions/>

LMSiv2和大语言模型框架的对应关系

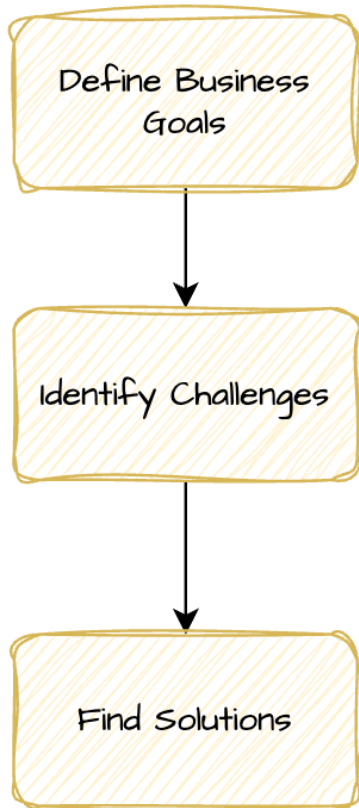
		AutoGen	LangChain	LlamaIndex	DSPY	HF Transformers	云厂商
1	应用层	😊	😊	😊			😊
2	优化层				😊		😊
3	控制层		😊	😊	😊		😊
4	提示约束层		😊	😊	😊		😊
5	提示层		😊	😊	😊		😊
6	神经网络层				😊	😊	😊
7	算力层						😊

02 案例分析

基于 Serverless 的大语言模型翻译应用



基于大语言模型的游戏内容翻译



业务声明

某游戏公司需要翻译游戏的多语言版本，因为存在较多游戏中特有的地名，人名，机器翻译的效果不好，主要是人工翻译为主，翻译的时间根据工作量的不同从数天到数周不等，业务团队希望利用大语言模型加速翻译过程并降低翻译成本。

```
"mapping" : {  
  "CHS" : "奇怪的渔人吐司",  
  "CHT" : "奇怪的漁人吐司",  
  "DE" : "Missslungene Fischerschnitte",  
  "EN" : "Suspicious Fisherman's Toast",
```



快速发现挑战和对应的解决方案

专有名词的翻译

长文本翻译

准确度提升

翻译风格

合规合法

部署

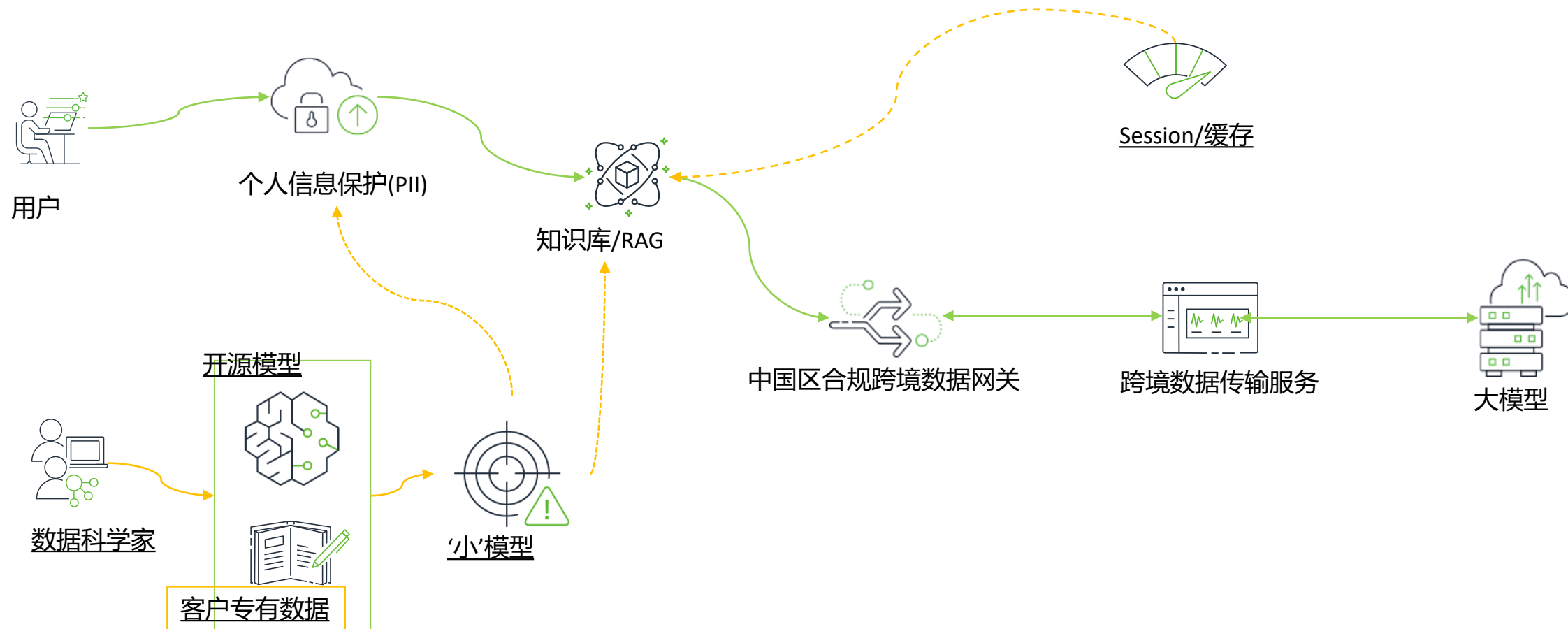
应用评估

效果反馈

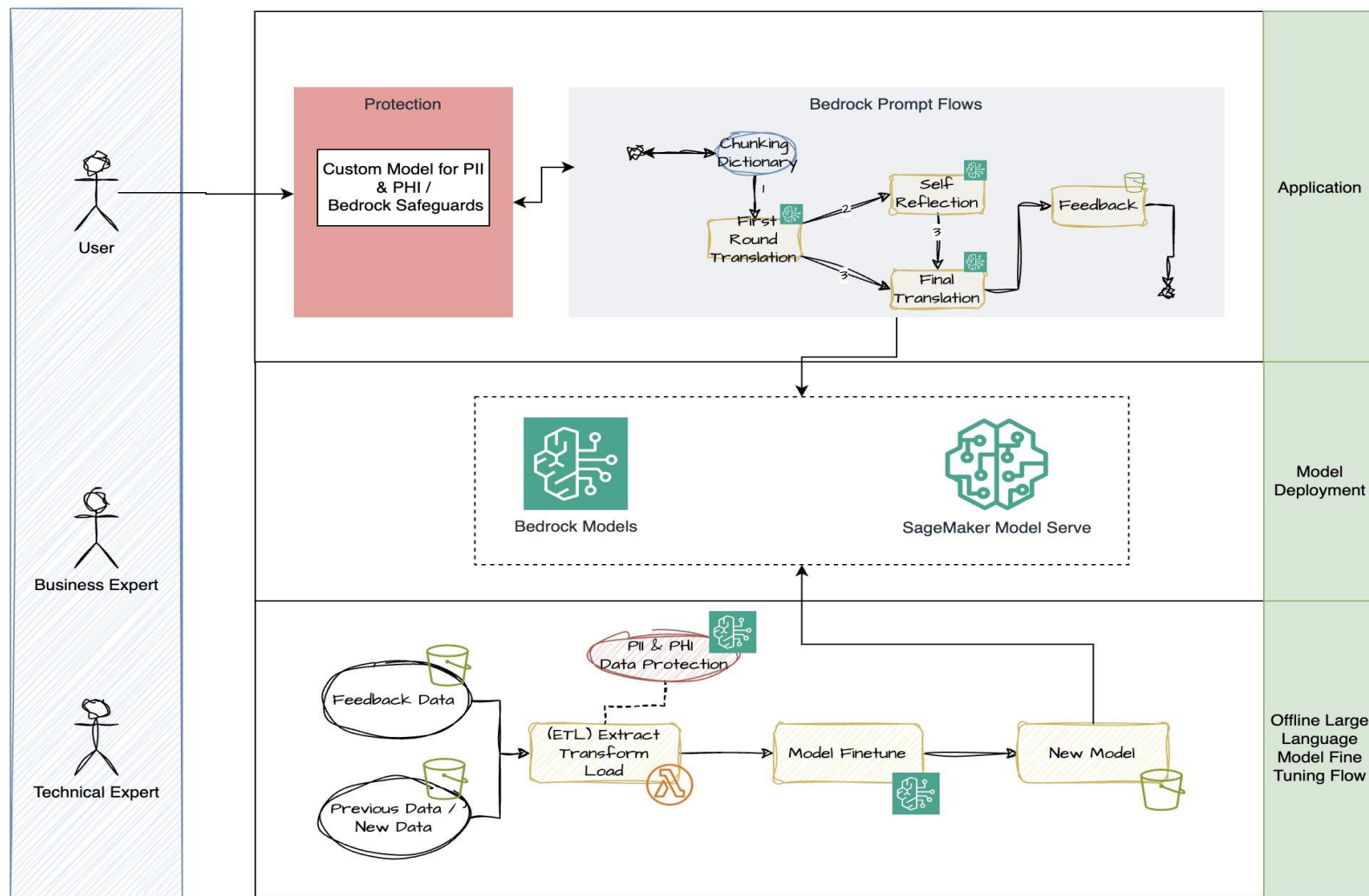
AI相关法规

名称	《互联网信息服务深度合成管理规定》	《生成式人工智能服务管理暂行办法》
信息安全要求	<p>第七条 深度合成服务提供者应当落实信息<u>安全主体责任</u>，建立健全用户注册、<u>算法机制机理审核、科技伦理审查、信息发布审核、数据安全、个人信息保护、反电信网络诈骗、应急处置等管理制度</u>，具有安全可控的<u>技术保障措施</u>。</p>	<p><u>输入 Input</u></p> <p>第十一条提供者对使用者的输入信息和使用记录应当依法履行保护义务，不得收集非必要个人信息，不得非法留存能够识别使用者身份的输入信息和使用记录，不得非法向他人提供使用者的输入信息和使用记录。</p> <p><u>输出 Output</u></p> <p>第九条提供者应当依法承担网络信息内容生产者责任，履行网络信息安全义务。 涉及个人信息的，依法承担个人信息处理者责任，履行个人信息保护义务。</p>

完整的大模型应用的复杂性



通用的大语言模型应用的解决方案



03 最佳实践

用Serverless承载大模型应用

以上内容仅为本文档的试下载部分，为可阅读页数的一半内容。如要下载或阅读全文，请访问：<https://d.book118.com/337113111166006163>