



算力网络泛在算力调度研究

周凡钦
北京邮电大学



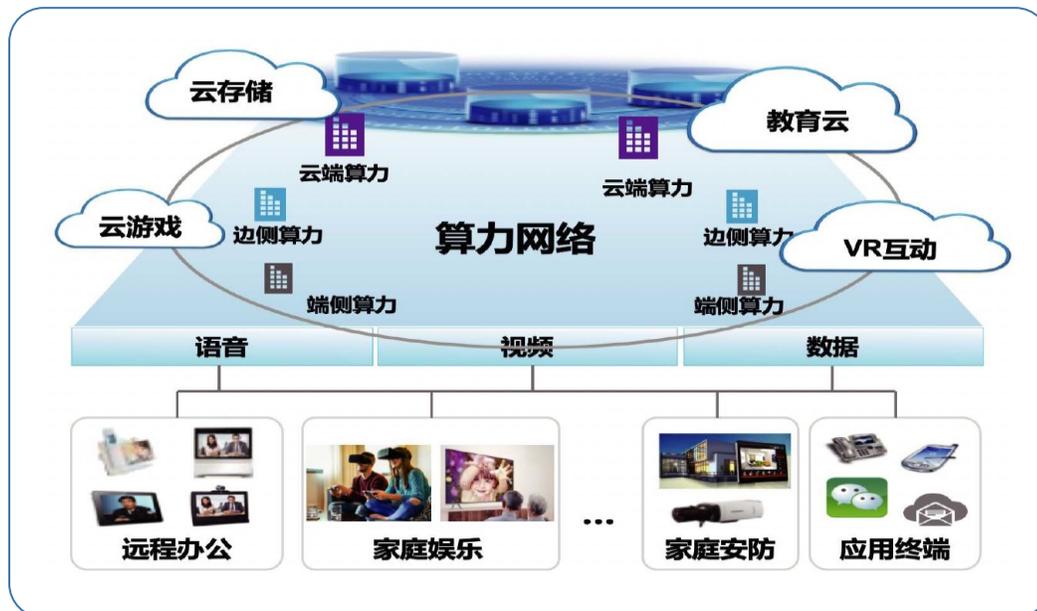
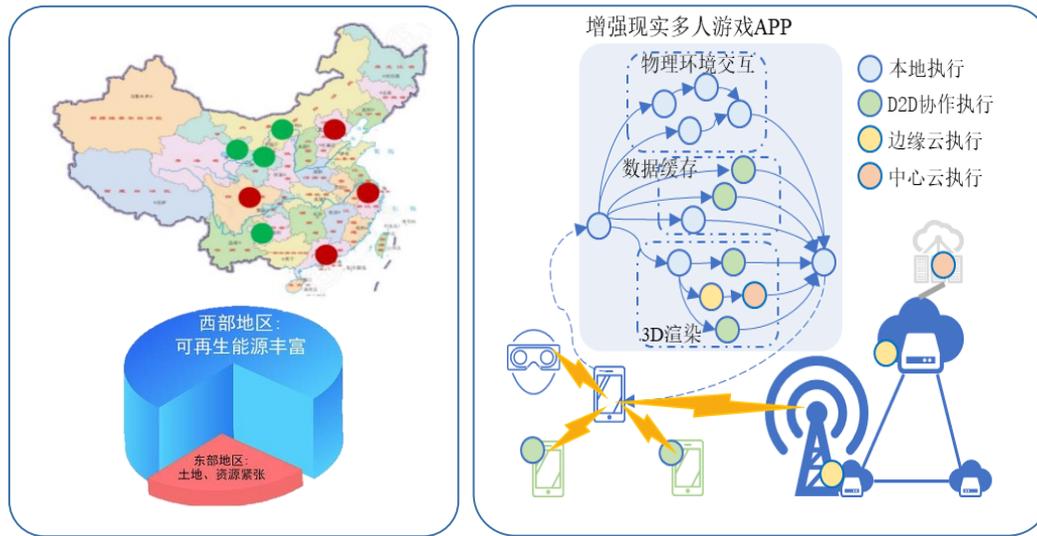
CONTENTS

1. 研究背景
2. 泛在算力调度技术研究—算力度量
3. 泛在算力调度技术研究—算力调度
4. 泛在算力调度整体功能框架

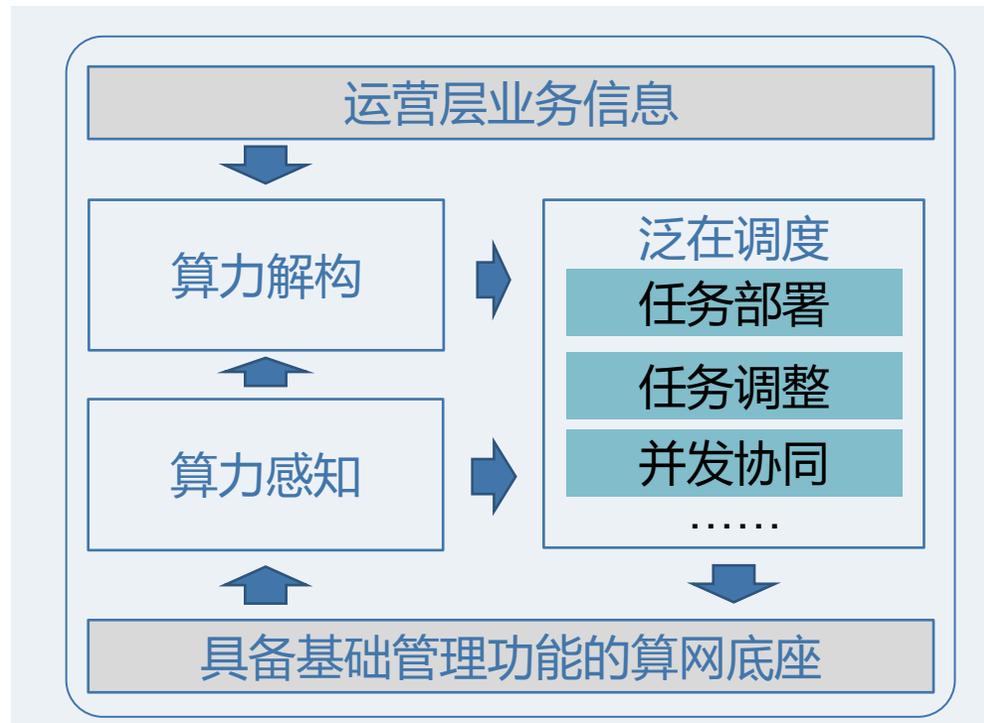
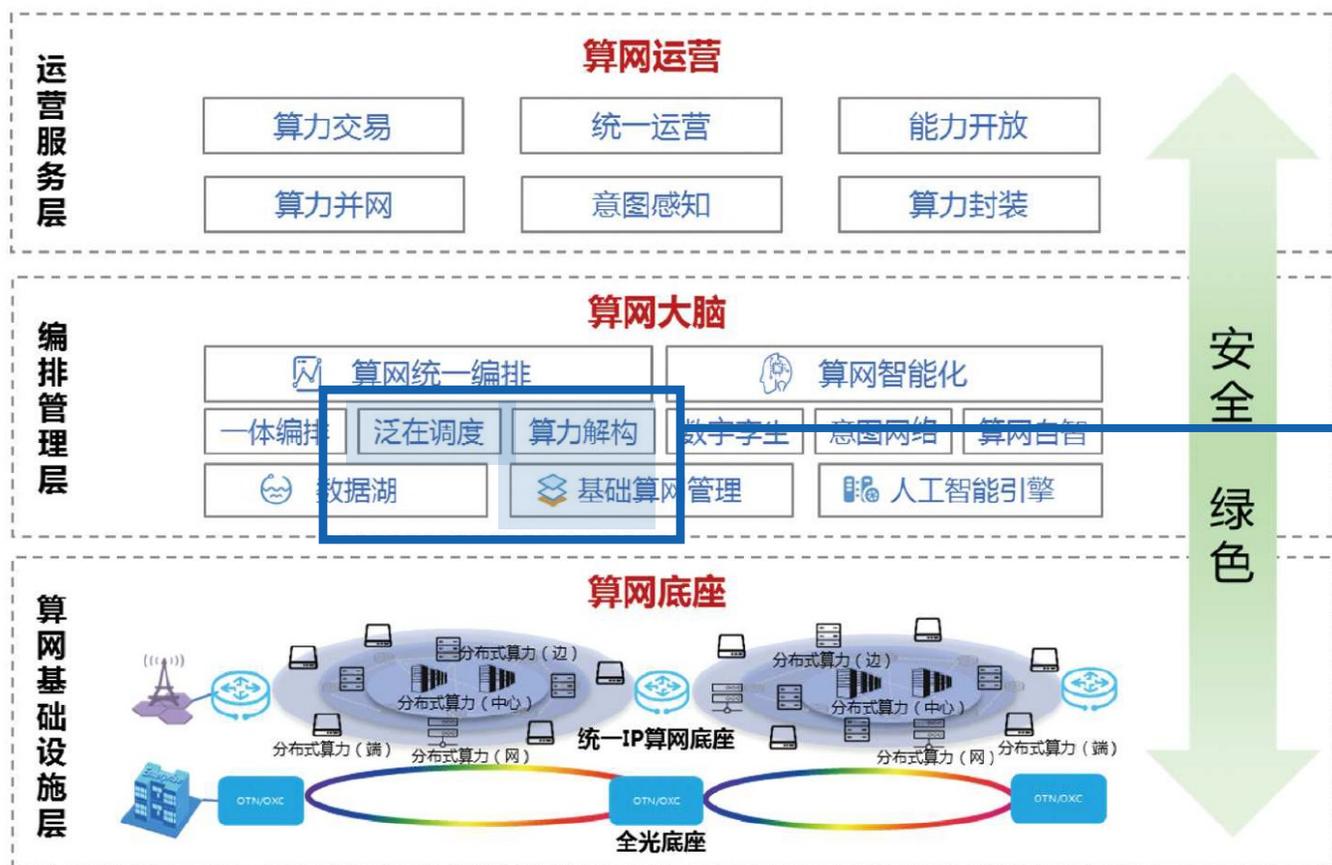


- **政策导向：**《“十四五”数字经济发展规划》要求：加快实施‘东数西算’工程，推进云网协同发展，提升数据中心跨网络、跨地域数据交互能力，加强面向特定场景的边缘计算能力，**强化算力统筹和智能调度”**。
- **业务需求：**新兴网络业务对“算”和“网”都提出了极高的要求，传统模式无法满足多种多样的新业务(如VR、车联网等、智能制造)实时处理、数据隐私（如工业用户数据不出园区，数据不出家庭网关等）等**多样化服务需求**。
- **技术因素：**边缘计算、微服务、云原生等技术快速发展，边缘+终端形成泛在分布的碎片化算力资源（不同规模、供方、异构）；微服务和云原生等组件化的应用架构，提升开发效率，为算力解构提供了基础。
- **社会效益：**用户端算力不断增强，社会面大量算力资源闲置，如何让闲置算力充分流动，实现云边端算力有效协同和资源效益的增长，是算力网络的重要愿景。

需要算网融合的新型网络技术体系——算力网络



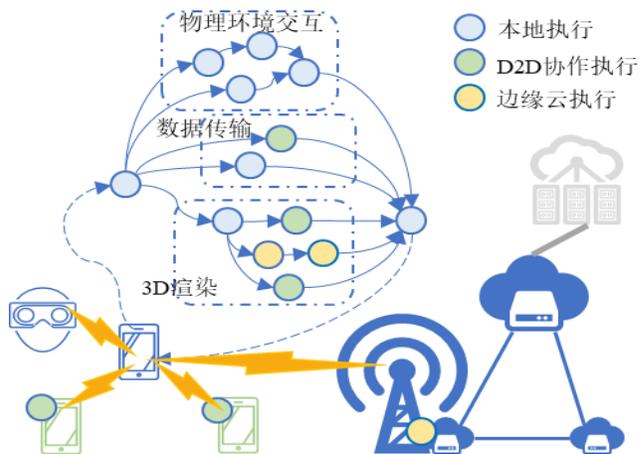
算力网络的逻辑架构由基础设施层、编排管理层、运营服务层构成。本团队主要关注算网大脑中泛在调度技术以及与之密切相关的算力结构和算力感知技术。



- 算力感知: 准确快速获取算和网的资源信息
- 算力解构: 复杂算力需求的**任务分解**
- 泛在调度: 面向具体业务和场景需求, 部署计算任务和调度算力资源

图▲ 算力网络体系架构

◆ 泛在计算的不同应用和场景及对算力的不同需求，同时应用的软件架构呈现组件化趋势，

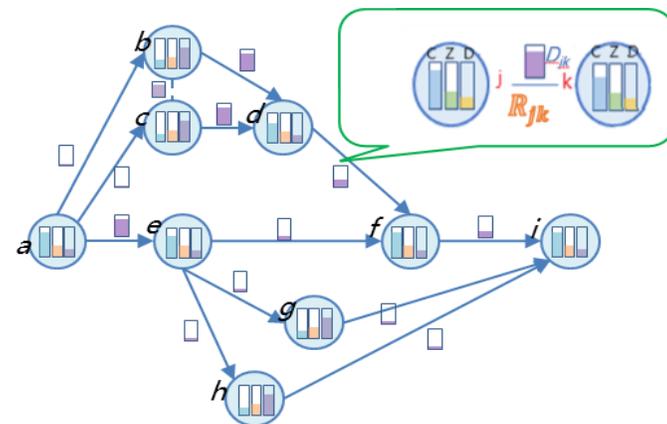


场景类别	示例	特点
云计算	政务云、企业云等云化服务;云端网络应用	计算密度低, 数据量不大, 时延要求不高; 或者计算量极高, 计算密度极大;
云-边协同	云游戏、智能安防	计算密度较高, 数据量大, 时延要求高;
云-端协同	物联智能业务、智能家居业务	海量泛在计算, 计算任务密度低, 数据量不大, 时延要求不高;
端-边协同	智能工厂	高实时, 高隐私和安全性
端-边-云协同	多人AR游戏、云制造	低时延, 多地现场计算任务协作, 各类计算量任务并存
端-端-边-边协同	自动驾驶	低时延, 高移动性, 各类计算量任务并存

◆ 面向算力网络的网络应用任务分解

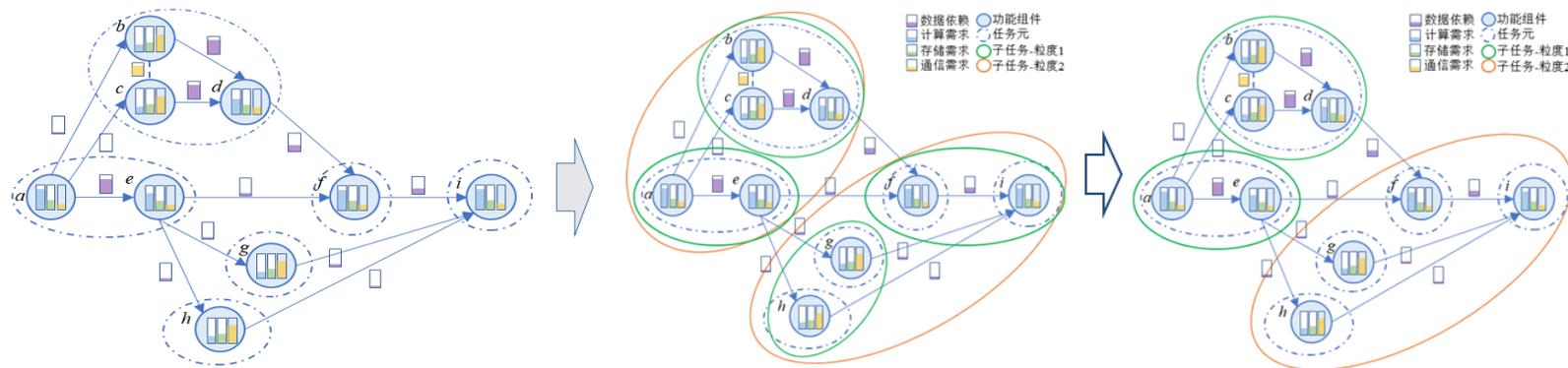
- 传统通信型应用: 具有单一子任务, 占用固定资源;
- 计算型应用: 数据节点持续推送数据到计算节点, 返回结果给目的节点, 形成循环结构;
- 内容型应用: 内容请求发送到内容索引服务节点, 内容源向内容请求终端发送内容数据;
- 复杂应用, 如右图增强现实多人游戏应用,

- 资源要求: 计算、通信和/或存储
- 性能需求: 带宽、信道质量、时延、算力
- 其他需求: 服务持续性、业务流粘性、资源节点统一性等



图▲ 抽取出的应用组件依赖关系图及资源需求

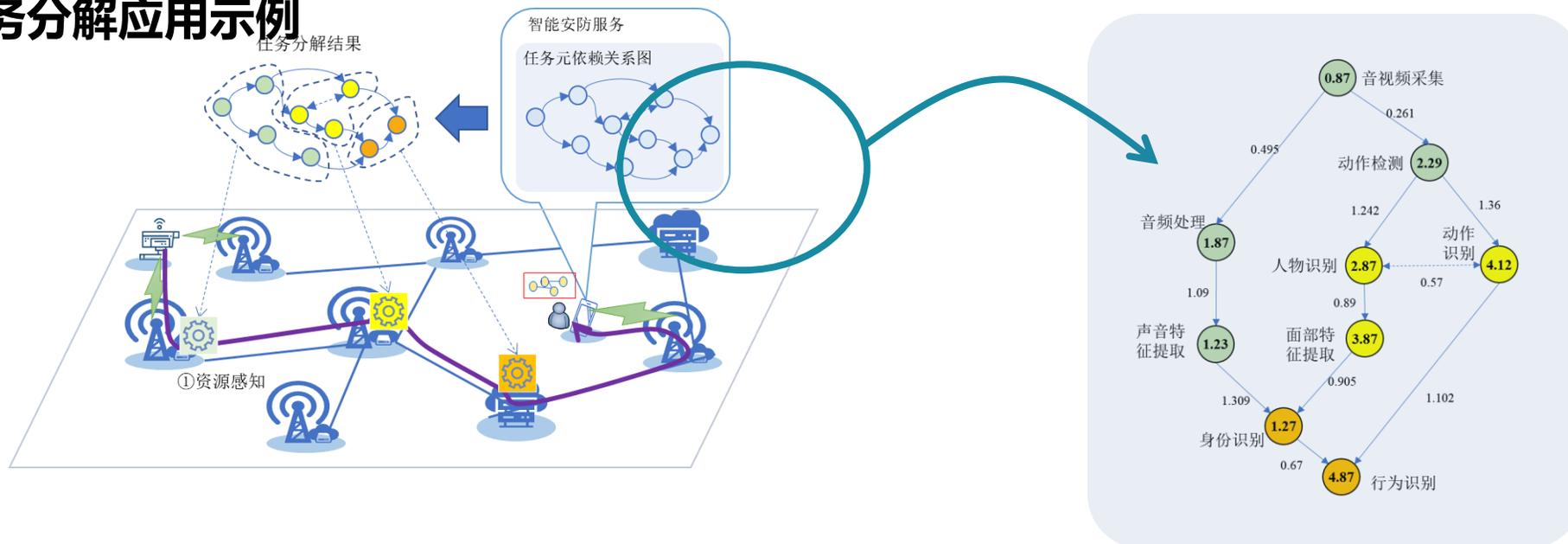
◆ 面向算力网络的网络应用多粒度任务分解



主要方案

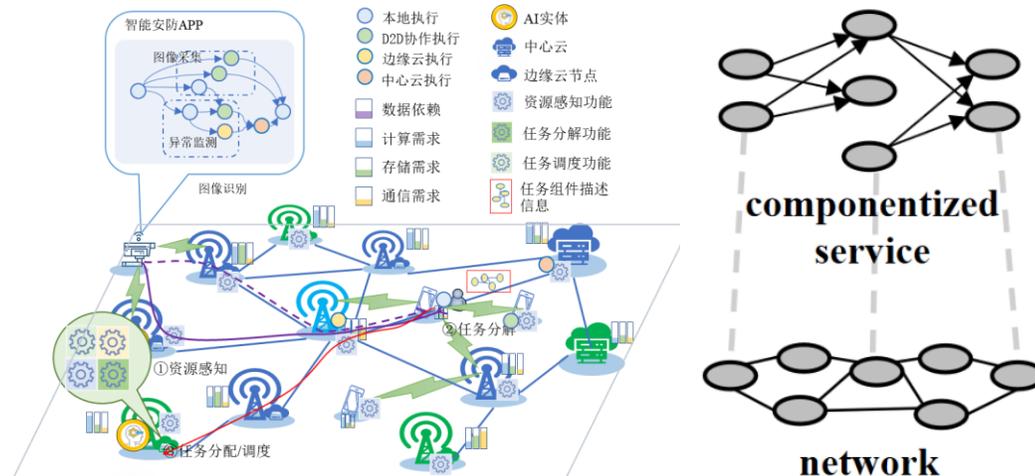
- 基于图聚类的任务分解
- 基于最优化模型的任务分级
- 面向多样化算力的多粒度分解
- 面向移动场景的任务分解
-

◆ 任务分解应用示例

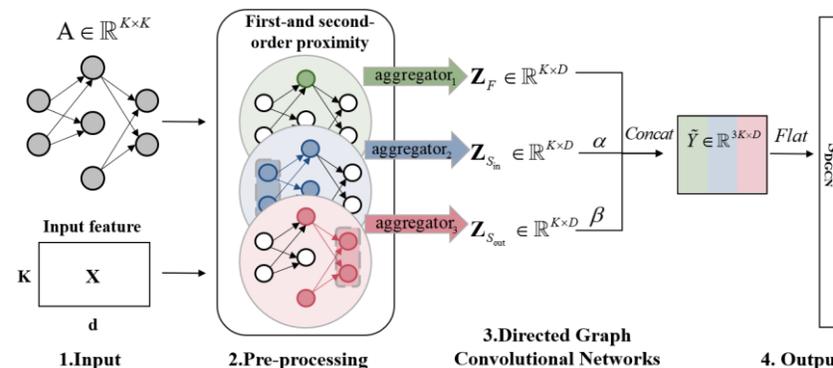


◆ 基于深度强化学习的面向组件化任务的资源调度算法

- 将资源利用均衡度作为优化目标，并同时让节点能效尽可能高，实现资源有效利用
- 建立泛在计算体系的资源调度模型
 - 计算、通信等多域异质资源的可用性度量方法；
 - 限制资源分配的关键约束分析和模型构建；
- 采用深度强化学习以实现快速推理，满足快速响应用户发起的泛在算力调度需求，以及支持调度策略自动演进。



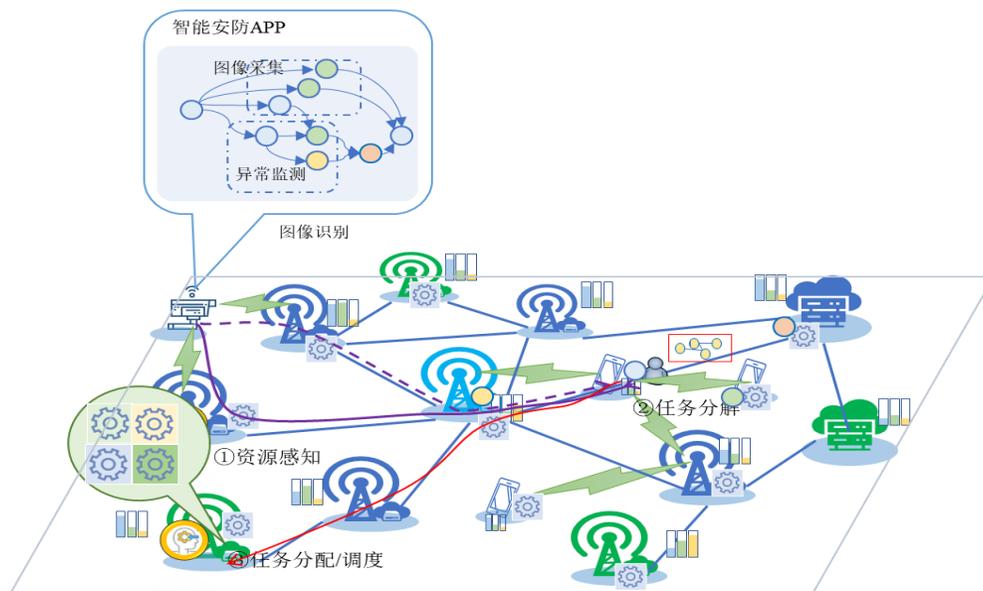
- 提出基于GNN增强的深度强化学习算法，提高收敛效率



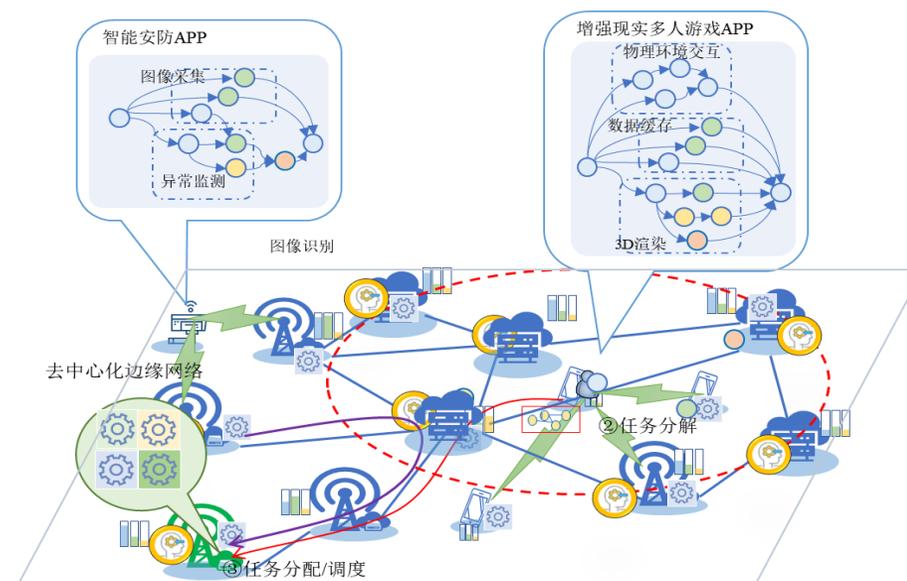
<p>所有任务请求的节点资源量不超过计算、存储资源总量</p>	$\max_{X_{UV}, P, w, c_V, z_V, b_E} \alpha_V \varphi_B(\eta_V) + \alpha_E \varphi_B(\eta_E) + \alpha_p \phi_E(P, w)$ $\text{s.t. } \begin{cases} x_{uV} C_u \leq c_{uV} \leq x_{uV} \mu_C(V) \\ x_{uV} S_u \leq z_{uV} \leq x_{uV} \mu_Z(V) \\ \sum_u x_{uV} c_{uV} \leq \mu_C(V) \\ \sum_u x_{uV} z_{uV} \leq \mu_Z(V) \\ \sum_{V_{u_k} = p} b_{V_{u_k} V} \leq \mu_D(V_u V) \end{cases}$	<p>节点的计算、存储资源满足某应用需求</p>
<p>任务数据传输时延不超过时间门限</p>	$T(U) \leq \hat{T}_U$ $L(u, \mathcal{D}_{du}(\cdot)) = 1$ $\tau_{u, V_p V} (c_{uV}, z_{uV}, b_{V_p V}) \leq \hat{\tau}_{d, u}$	<p>任务执行时间不超过执行时间门限</p>
<p>任务节点可靠性满足一定门限条件</p>	$\psi_{V_p V} (v_{V_p V}, p_s(V_p, V)) \geq \hat{\psi}$	<p>任务节点可靠性满足一定门限条件</p>

◆ 面向特定场景的泛在算力资源调度方案研究

- 适配算力资源条件变化的业务资源调整
 - 节点和链路资源变化，影响业务，调整算力资源
 - 满足业务资源需求
 - 降低对现有部署的影响



- 分布式多智能管控代理协同下的算力调度
 - 管控代理控制特定区域资源
 - 多代理协同为跨域业务提供所需资源
 - 解决多代理情况下资源冲突与资源效率的矛盾





算力网络泛在调度原型-工作原理



原型系统实现了算力注册、算力感知、算力分解、任务部署及动态调整，验证了算力网络泛在调度的基本功能。

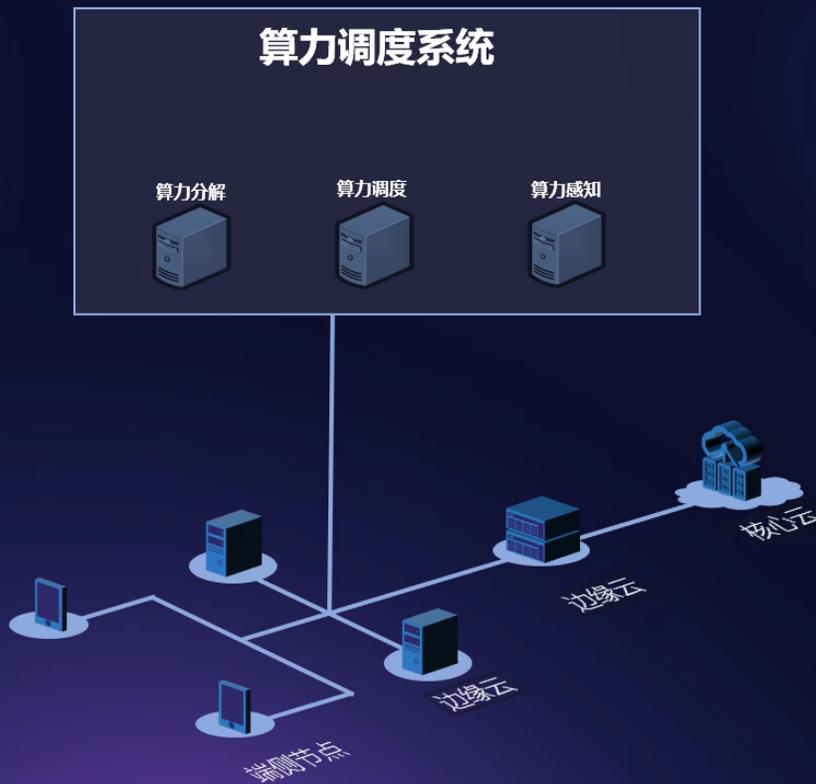
算力注册过程

算力感知过程

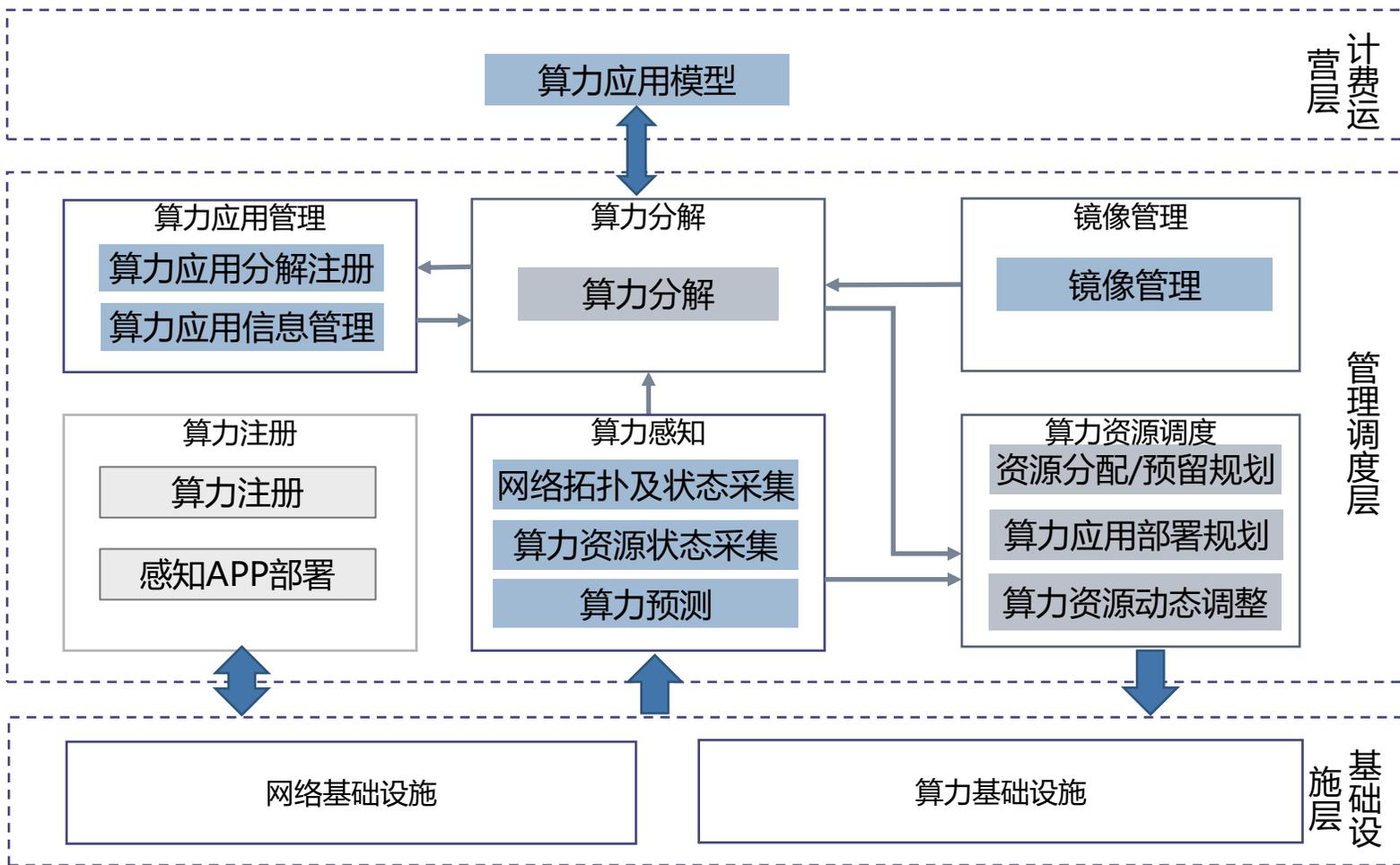
任务分解过程

任务部署过程

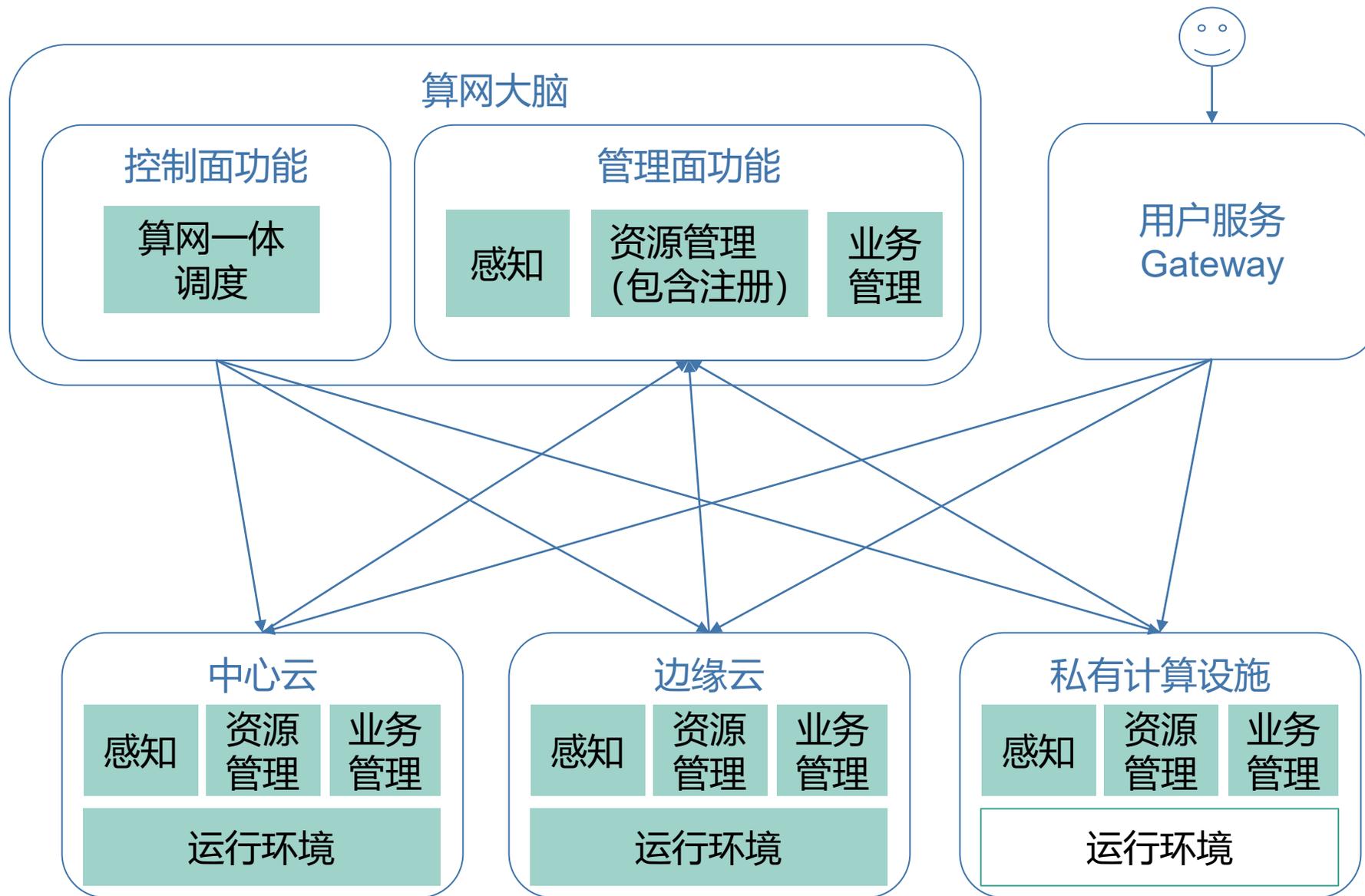
运行调整过程



原型系统功能框架设计



- **算力注册**: 新节点加入泛在算力体系, 向算力调度系统申请注册, 注册成功后, 算力调度系统向该节点发送管理代理。
- **算力感知**: 管理代理中的感知代理定期向算力调度系统上报实时算力信息。
- **任务分解**: 用户发送算力应用请求, 算力调度系统收到应用信息对应用进行分解, 并封装为子任务。
- **任务部署**: 算力调度系统根据应用分解后的子任务信息和用户周边的算力资源信息, 生成任务部署方案, 并完成部署。
- **运行调整**: 应用运行过程中, 资源条件发生改变, 算力调度系统会调整任务部署和算力资源分配, 保障应用平稳运行。





感谢聆听

合作·创新

北京邮电大学计算机院网络管理研究中心

Email: fqzhou2012@bupt.edu.cn

致谢：感谢文中所引用素材的机构及作者