

- 资源扩容方案 - 快速对比表
  - 📊 资源配置对比
  - 💰 成本对比
  - 🎯 性能对比
  - ✅ 适用场景
  - 🎯 推荐决策树
  - 🚀 配置文件
  - 💡 关键优化点
    - 方案A（高配）核心优化 ⭐
    - 方案B（超高配）额外优化 🔥
  - 📋 快速决策建议
    - 如果您是...
  - ⚠️ 重要提示
  - 🎯 最终建议
    - 立即行动：
    - 为什么先选方案A？

# 资源扩容方案 - 快速对比表

3种方案一目了然对比



## 资源配置对比

容器组件	当前低配	方案A 高配 ⭐	方案B 超高配
Server CPU	4核 (2核预留)	8核 (4核预留)	16核 (8核预留)
Server 内存	8GB (4GB预留)	16GB (8GB预留)	32GB (16GB预留)
Web CPU	2核 (1核预留)	4核 (2核预留)	6核 (3核预留)
Web 内存	4GB (2GB预留)	8GB (4GB预留)	12GB (6GB预留)
MySQL CPU	2核 (1核预留)	4核 (2核预留)	6核 (3核预留)
MySQL 内存	2GB (1GB预留)	4GB (2GB预留)	8GB (4GB预留)
MySQL 缓冲池	512MB	2GB ⬆️ 4x	4GB ⬆️ 8x

容器组件	当前低配	方案A 高配 <span>★</span>	方案B 超高配
Redis CPU	1核 (0.5核预留)	2核 (1核预留)	4核 (2核预留)
Redis 内存	1GB (512MB预留)	2GB (1GB预留)	4GB (2GB预留)
Redis 缓存	512MB	1.5GB <span>↑</span> 3x	3GB <span>↑</span> 6x
总CPU	9核	18核 <span>↑</span> 100%	32核 <span>↑</span> 255%
总内存	15GB	30GB <span>↑</span> 100%	56GB <span>↑</span> 273%



## 成本对比

项目	当前低配	方案A	方案B
推荐服务器配置	16核32GB	20核36GB	36核64GB
阿里云/腾讯云 月付	¥800-1200	¥1600-2400	¥3000-4500
年付价格	¥8000-12000	¥16000-24000	¥30000-45000
相对成本	基准	+¥800-1200/月	+¥2200-3300/月



## 性能对比

指标	当前低配	方案A 高配	方案B 超高配
并发用户数	5-10人	10-20人	20-50人
LLM推理延迟	2-5秒	1-2秒 <span>↓</span> 60%	<1秒 <span>↓</span> 80%
API响应时间	500-1000ms	200-500ms <span>↓</span> 50%	<200ms <span>↓</span> 70%
数据库查询	100-500ms	50-100ms <span>↓</span> 50%	<50ms <span>↓</span> 80%
缓存命中率	80-85%	90-95%	95-99%
卡顿频率	较频繁	偶尔	几乎无

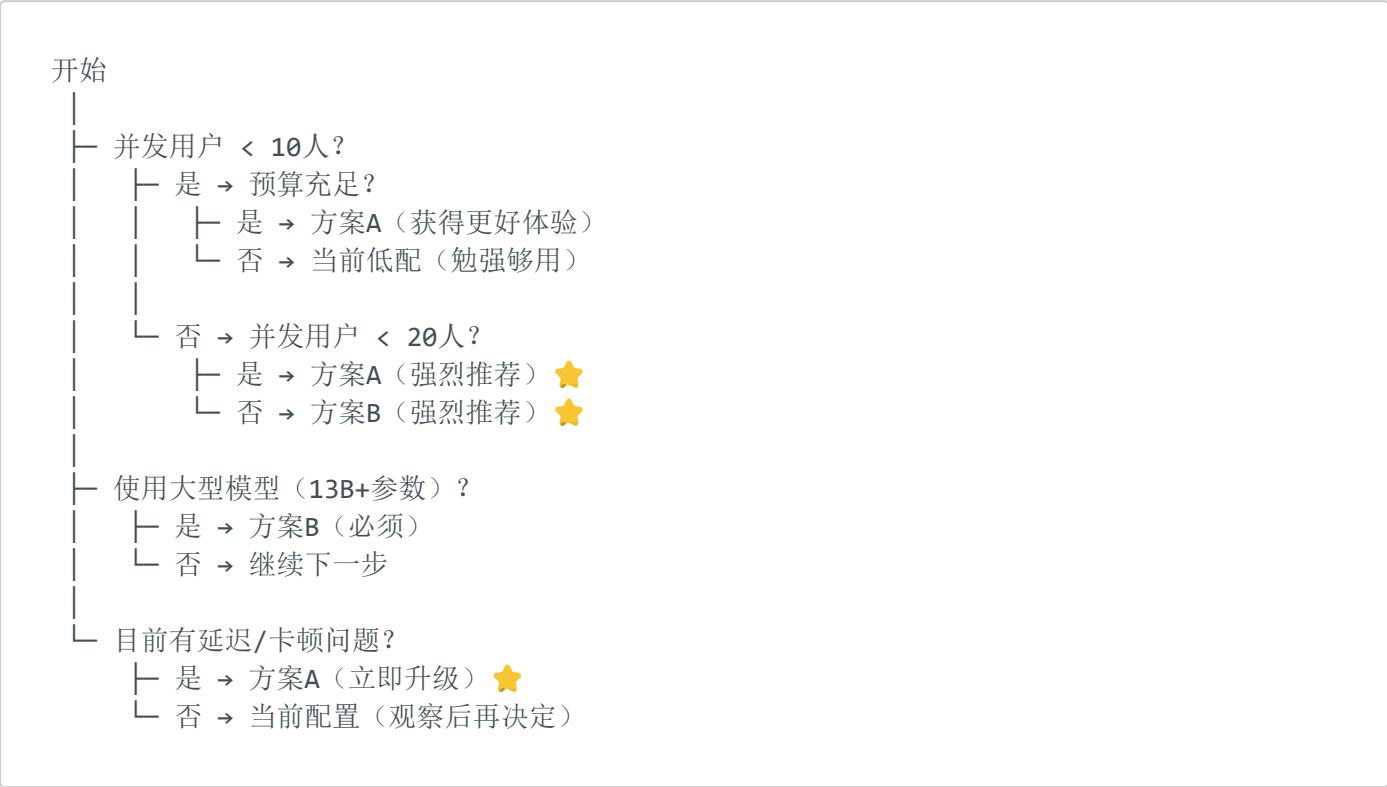


## 适用场景

场景	当前低配	方案A	方案B
测试/开发环境	✅ 可用	✅✅ 更好	⚠️ 过高
小型团队 (<10人)	⚠️ 勉强	✅✅ 推荐	⚠️ 过高
中型企业 (10-20人)	❌ 不足	✅✅ 推荐	✅ 可选
大型企业 (20-50人)	❌ 不足	⚠️ 勉强	✅✅ 推荐
超大规模 (50+人)	❌ 不可用	❌ 不足	⚠️ 勉强
使用大型模型 (13B+)	❌ 不可用	⚠️ 勉强	✅✅ 推荐



## 推荐决策树



## 配置文件

方案	配置文件	使用方法
当前低配	<code>docker-compose-offline.yml</code>	现有配置
方案A 高配	<code>docker-compose-high.yml</code>	✅ 已创建
方案B 超高配	<code>docker-compose-ultra.yml</code>	✅ 已创建

## 使用方法:

```
# 停止现有服务
docker-compose -f docker-compose-offline.yml down

# 启动高配方案
docker-compose -f docker-compose-high.yml up -d

# 或启动超高配方案
docker-compose -f docker-compose-ultra.yml up -d
```



## 关键优化点

### 方案A（高配）核心优化 ⭐

#### 1. Server容器

- CPU: 4核 → 8核 (⬆ 100%)
- 内存: 8GB → 16GB (⬆ 100%)
- 效果: LLM推理速度提升一倍

#### 2. MySQL数据库

- InnoDB缓冲池: 512MB → 2GB (⬆ 4x)
- 连接数: 1000 → 2000
- 效果: 查询延迟降低50%

#### 3. Redis缓存

- 缓存大小: 512MB → 1.5GB (⬆ 3x)
- 效果: 缓存命中率提升到95%

### 方案B（超高配）额外优化 🔥

#### 1. Server容器

- CPU: 8核 → 16核 (再 ⬆ 100%)
- 内存: 16GB → 32GB (再 ⬆ 100%)

- 支持13B+参数的大型模型

## 2. MySQL数据库

- InnoDB缓冲池: 2GB → 4GB
- IO线程: 4 → 8（读写各8个）
- 连接数: 2000 → 5000

## 3. Redis缓存

- 缓存大小: 1.5GB → 3GB
- IO多线程: 开启4线程
- 连接数: 10000 → 20000



# 快速决策建议

## 如果您是...

场景	推荐方案	理由
测试/开发团队	方案A	性价比最高，体验明显改善
小型企业（<20人）	方案A	满足需求，成本可控
中型企业（20-50人）	方案B	支持更多并发，无延迟
使用大型模型	方案B	必须，否则无法运行
预算有限	方案A	先升级到方案A，不够再升B
追求极致性能	方案B	几乎无延迟，最佳体验



## 重要提示

### 1. 服务器选择建议

- 方案A: 申请 **20核36GB**（预留Buffer）
- 方案B: 申请 **36核64GB**（预留Buffer）
- 磁盘: 至少200GB SSD（数据量大选500GB+）

## 2. 网络带宽

- 10人以下: 10Mbps
- 20人: 20Mbps
- 50人: 50Mbps+

## 3. 升级步骤

- ☒ 先备份数据
- ☒ 停止现有服务
- ☒ 使用新配置文件启动
- ☒ 监控资源使用
- ☒ 性能测试验证



## 最终建议

### 立即行动：

1. ☒ 申请方案A资源（20核36GB）
2. ☒ 使用 `docker-compose-high.yml`
3. ☒ 部署并观察1-2周
4. ☒ 如仍不足，再升级到方案B

### 为什么先选方案A？

- ☒ 性价比最高（2倍资源，>2倍性能）
- ☒ 能解决80%的性能问题
- ☒ 成本增加可控（月增¥800-1200）
- ☒ 后续可平滑升级到方案B

---

准备好申请资源了吗？ 