



大数据的理解与 分布式进化计算方法

张军

中山大学 超级计算机学院

2014年10月

大数据是什么？



大数据 = 海量数据？



大数据是海量数据的另一种说法

大数据只是体量比海量数据更大

大数据挖掘就是海量数据挖掘



我们的理解

大数据的认识

从计算机的发展历程说起



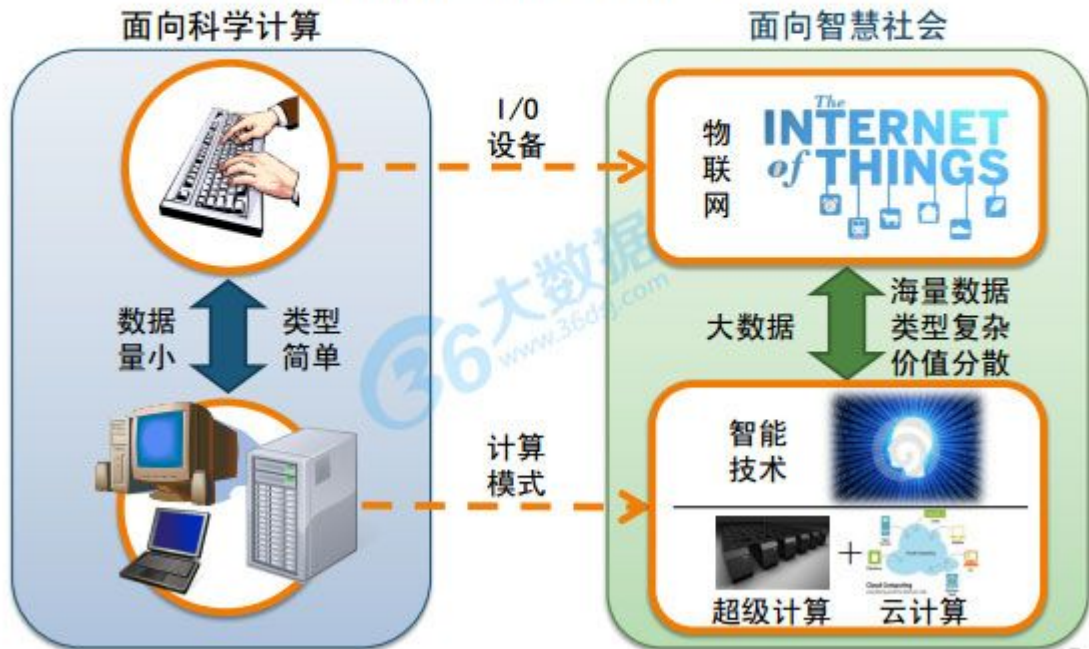
从计算机的发展历程说起



大数据的认识



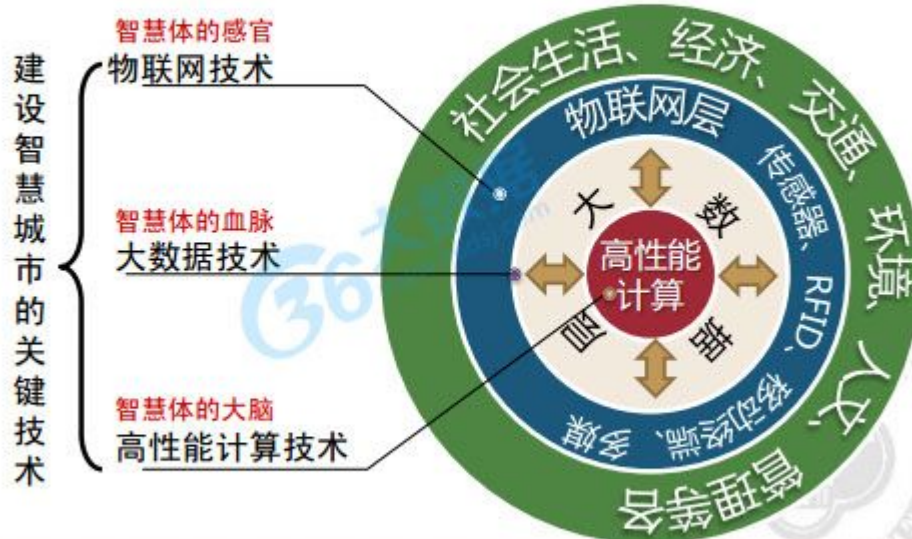
大数据是物联网与新型计算模式发展的产物



物联网 - 大数据 - 分布式计算



大数据与智慧城市



“智慧城市”是全球城市化建设的目标，也是我国跨越式发展的重要途径

大数据与海量数据的区别



类型

海量数据：数据类型简单，以结构化数据为主

大数据：数据类型复杂，半结构化和非结构化数据占主导地位

模型

海量数据：具有统计规律，能够通过数学模型进行描述

大数据：缺少统计规律，难以用数学模型来描述

方法

海量数据：经过长年探索已经形成一套可行的处理方法

大数据：尚且缺少行之有效的处理方法，亟待发展新方法

目标

海量数据：有明确的挖掘目标，关注解释事物之间的因果关系

大数据：没有具体的挖掘目标，关注点从因果关系转向关联关系；其价值在于能够发现超出预想的知识，填补空白。

如何解决大数据问题？



大数据的特点

体量巨大

类型复杂

价值丰富但分布不均

缺乏数学模型

具有容错性，以找到可接受解为目标

基于数学逻辑的人工智能方法难以解决缺乏数学模型的大数据问题

解决

大数据的特点使传统数据分析方法不再适用

新方法的可行途径

人工智能领域的计算智能方法将成为解决大数据问题的主要途径

进化计算

群体智能

深度学习

海量数据与大数据分析的异同

海量数据分析

有明确的分析目标

注重获取因果关系

追求一个精确的结果

VS

大数据分析

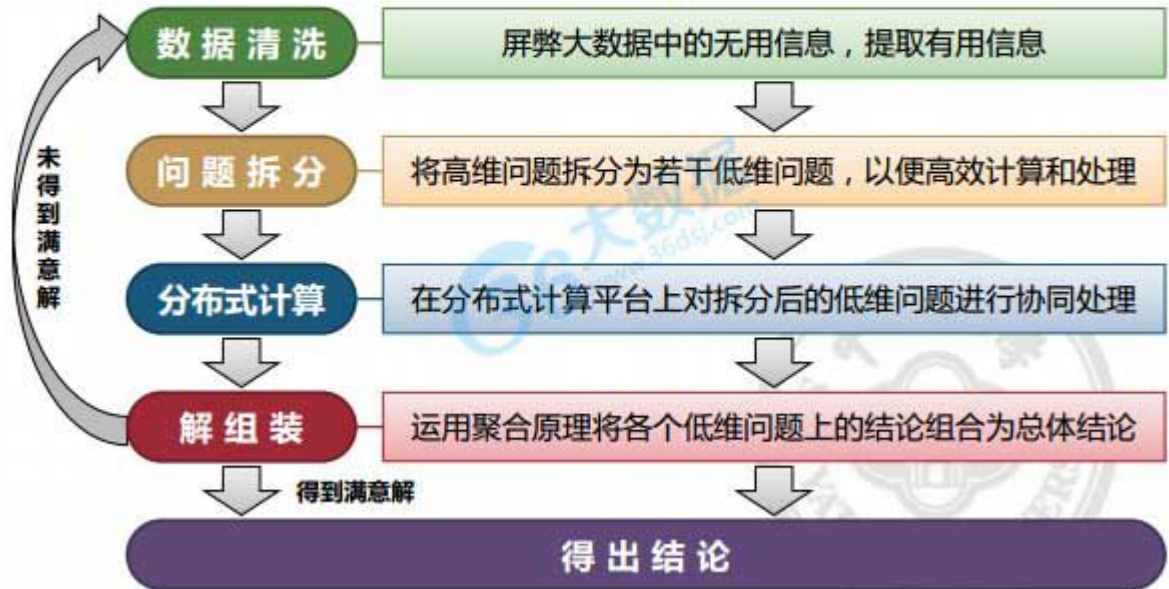
没有具体的分析目标

不仅挖掘因果关系，
更注重发现关联关系

发现未知的知识

期望快速找到可接受解

整体思路





问题拆分方法

值得研究的拆分方法

随机拆分

按任务拆分

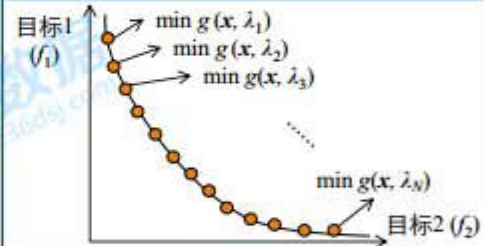
按目标拆分

⋮

复杂多目标问题

$$\min g(x, \lambda) = \lambda_1^{(1)} f_1(x) + \lambda_1^{(2)} f_2(x)$$

where $\lambda_1^{(1)} + \lambda_1^{(2)} = 1$



分拆

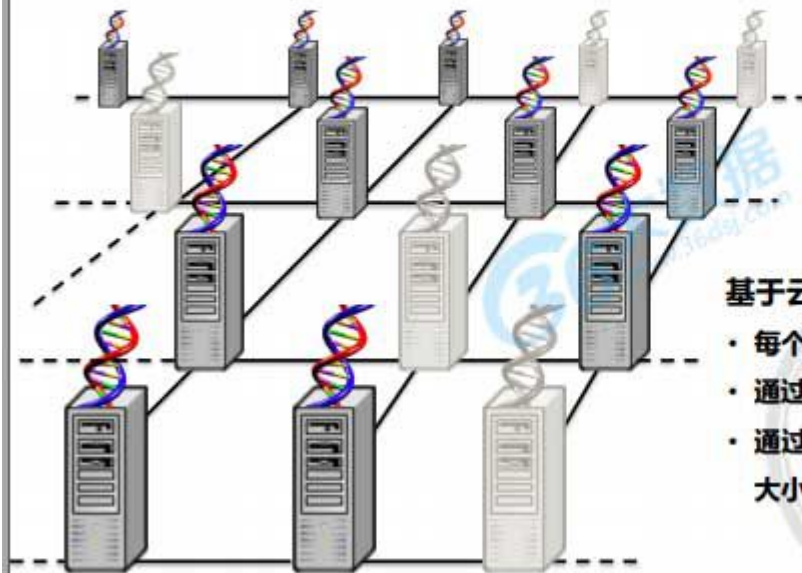
多个单目标问题

$\min g(x, \lambda_1)$ \longrightarrow 子系统1

\vdots

$\min g(x, \lambda_N)$ \longrightarrow 子系统N

分布的、自组织的智能计算方法



基于云计算平台的分布式进化算法

- 每个计算结点保存一个染色体
- 通过通信网络实现染色体的交叉
- 通过唤醒或删除计算结点来实现种群大小的自适应调整

研究背景与意义



国家需求

- “智慧城市”的建设需求
- 多个领域的**智能化需求**



智能交通控制



智能数据处理



智能传感器网络

智能化需求中的优化问题

- 缺乏**精确数学模型**
- 具有**高维、非线性**等特点
- 传统方法难以求解

需求

进化计算

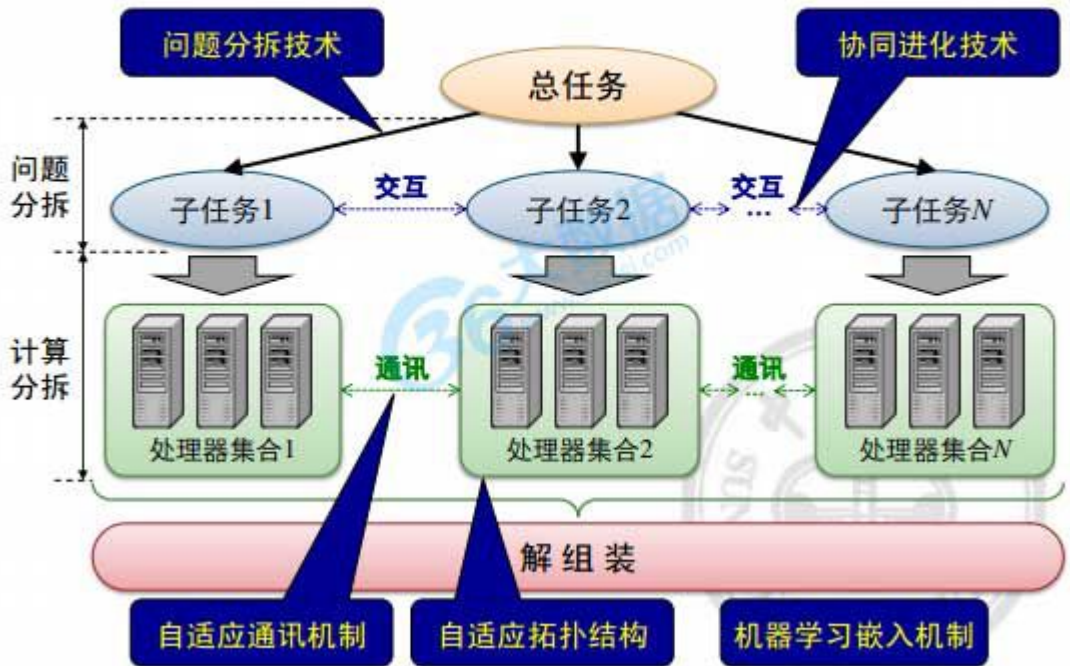
- 模拟自然进化的智能方法
- 不依赖问题的数学特性
- 成为解决复杂优化问题的重要途径

• 传统进化计算方法求解**大规模复杂优化问题**时存在**性能瓶颈**

• 传统进化计算方法以**串行执行**为主，无法发挥**分布式计算平台**优势

分布式进化计算

主要研究内容

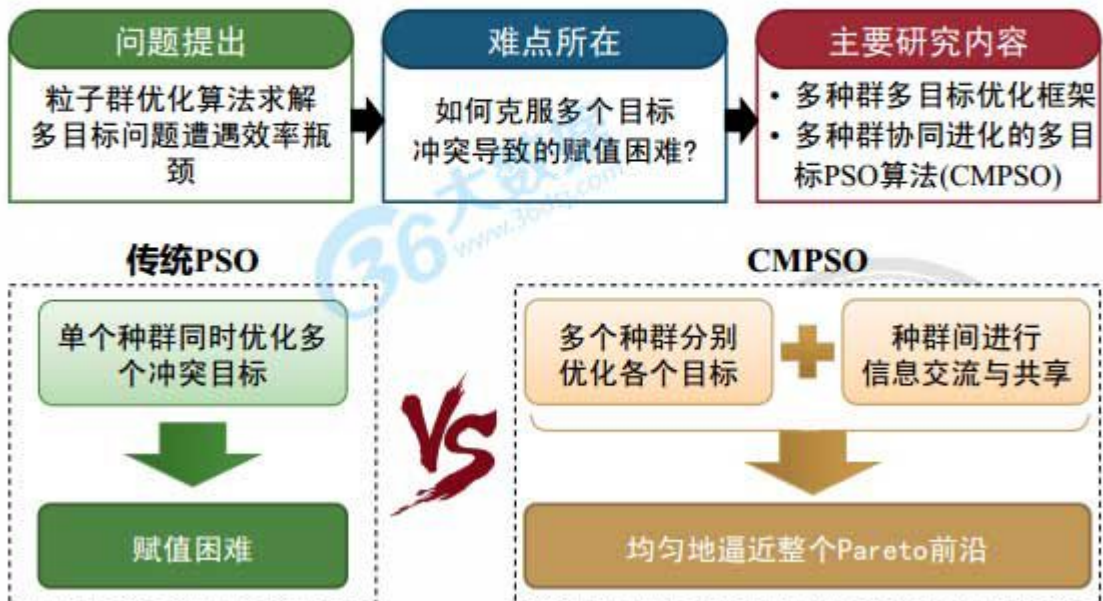


28

研究1：多种群协同进化的多目标PSO



研究背景与意义



29



算法思想



主要创新点：

- ☑ 每个种群优化一个目标：**解决了赋值困难，充分探索了目标空间**
- ☑ 种群通过全局Archive共享信息：**促使解均匀覆盖整个Pareto前沿**
- ☑ 上述思想可兼容各类进化算法：**多种群协同优化多目标问题的通用算法框架**



单个种群的进化策略

- ◎ 每个种群优化一个目标，进化过程与单目标PSO算法类似。
- ◎ 每个种群中粒子的位置更新方法：

新位置 = 原位置 + 新速度

每个粒子从全局Archive中**随机**选择一个学习对象：

- ☑ 计算量小
- ☑ 保持多样性

原速度的惯性影响

个体历史最佳位置的引导

种群历史最佳位置的引导

全局Archive中非劣解的引导



全局Archive的更新策略 (1)





全局Archive的更新策略 (2)

◎ 精英学习策略：

新解 = 原解 + 随机扰动向量

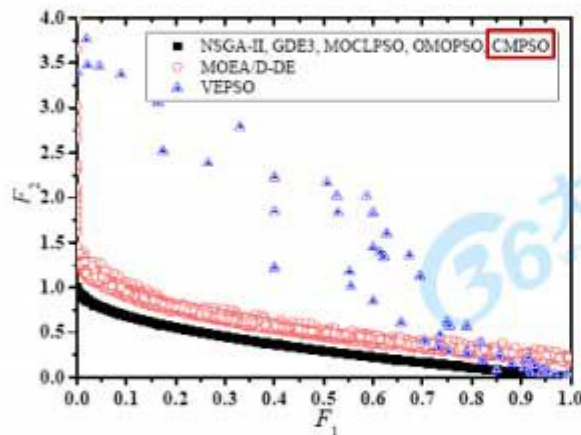
随机扰动全局Archive中的解，帮助算法跳出局部最优。

◎ 非劣解选择策略：

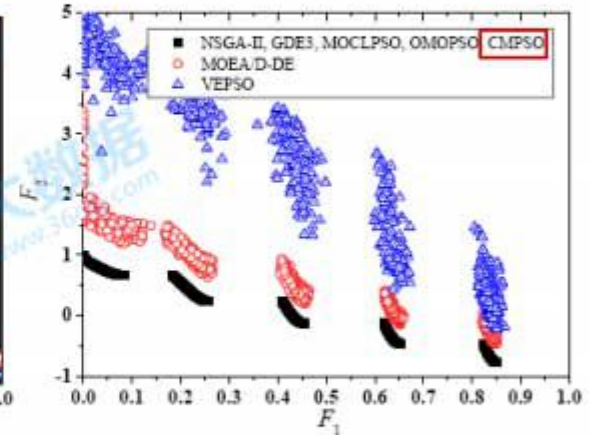
- ① 运用非支配排序（non-dominated sorting）选出非劣解；
- ② 如果非劣解的数目超过Archive的容量，根据拥挤距离（crowding-distance）删去分布密度最大的若干解。



实验结果与讨论



(a) ZDT1

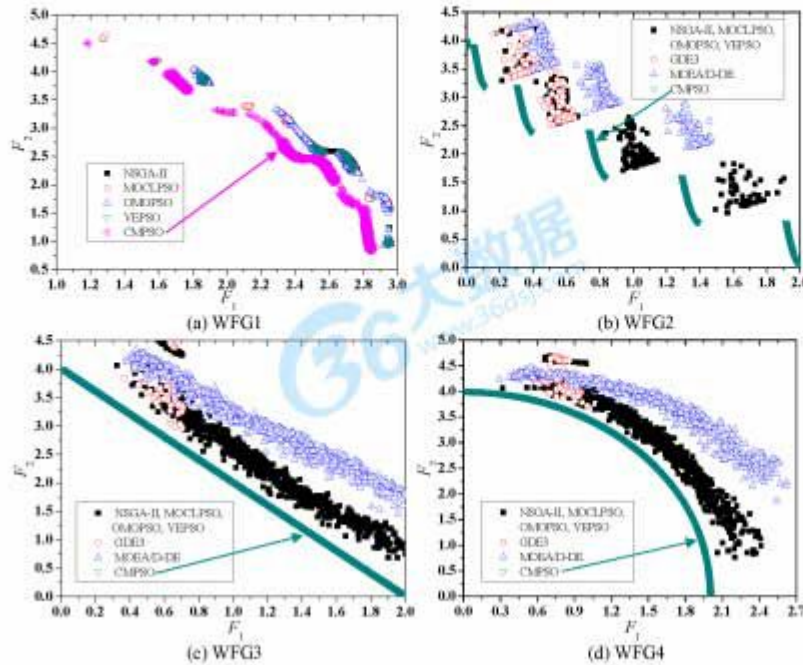


(b) ZDT3

CMPSO在ZDT系列问题上能够很好逼近的Pareto前沿



实验结果与讨论

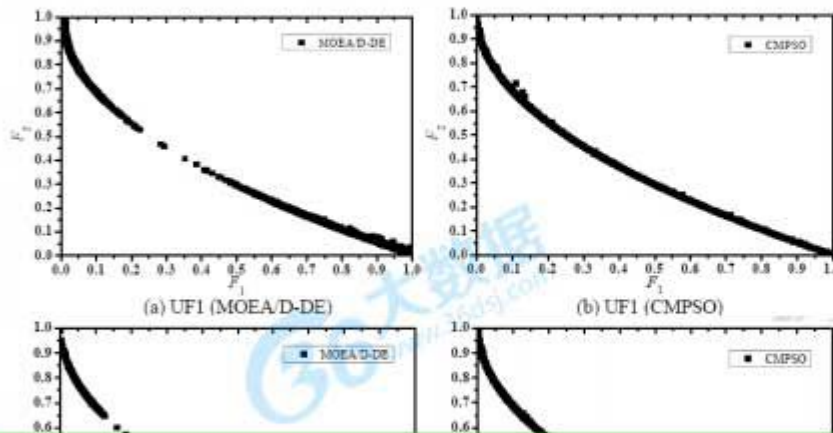


CMPSO 在 WFG系列问题上的全局收敛能力比其他算法好

研究1：多种群协同进化的多目标PSO



实验结果与讨论



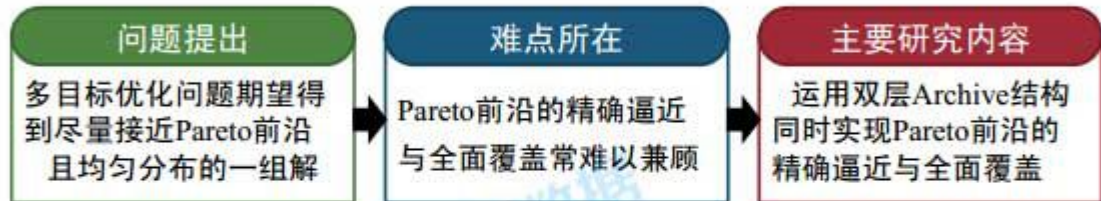
CMPSO 在 UF1 和 UF2 上对 Pareto 前沿的覆盖比 MOEA/D-DE 更均匀

研究成果已发表于国际期刊 *IEEE Transactions on Cybernetics*

Z.-H. Zhan, J. Li, J. Cao, J. Zhang, H. Chung, and Y. H. Shi, "Multiple populations for multiple objectives: A coevolutionary technique for solving multiobjective optimization problems," *IEEE Transactions on Cybernetics*, vol. 43, no. 2, pp. 445 – 463, Apr. 2013.



研究背景与意义

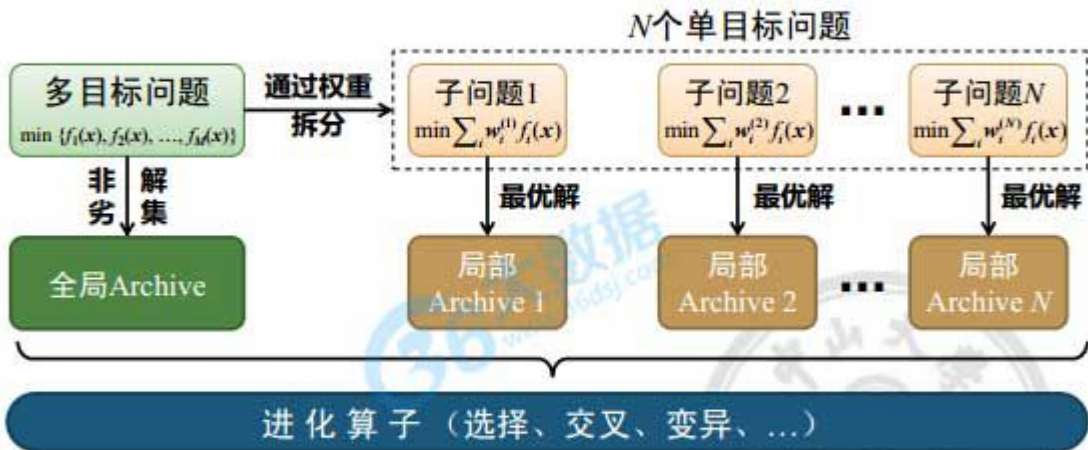


现有多目标进化算法





算法思想



主要创新点：

- ☑ 在多目标问题和单目标子问题两个层面上同时进化：兼顾全面性和搜索精度
- ☑ 上述算法框架可兼容多种进化算子：灵活性强，可根据实际问题进行拓展

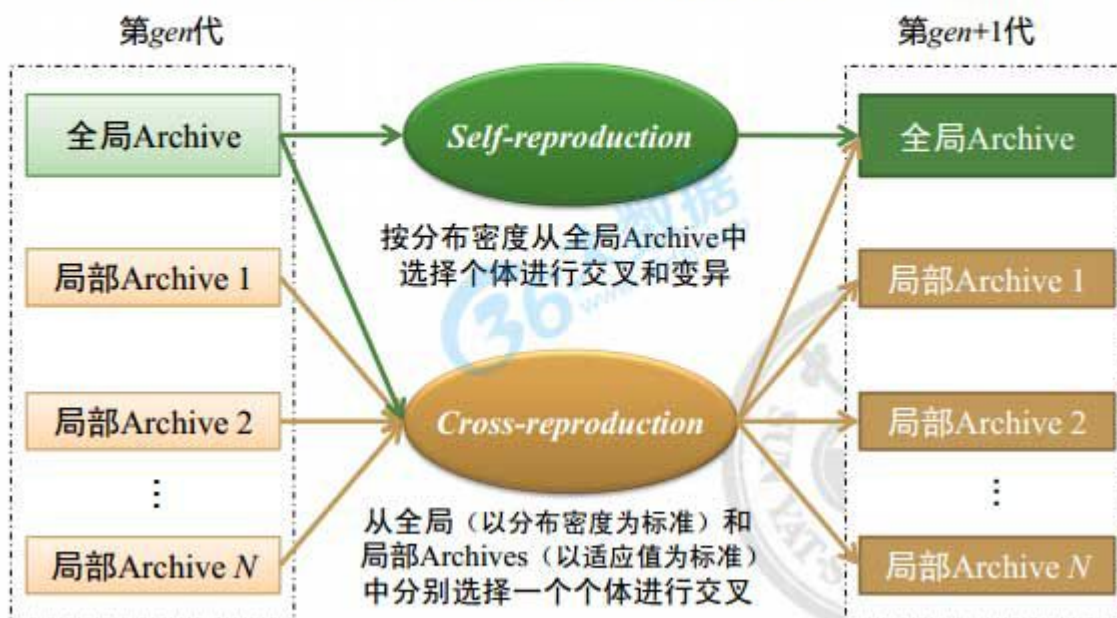


算法流程

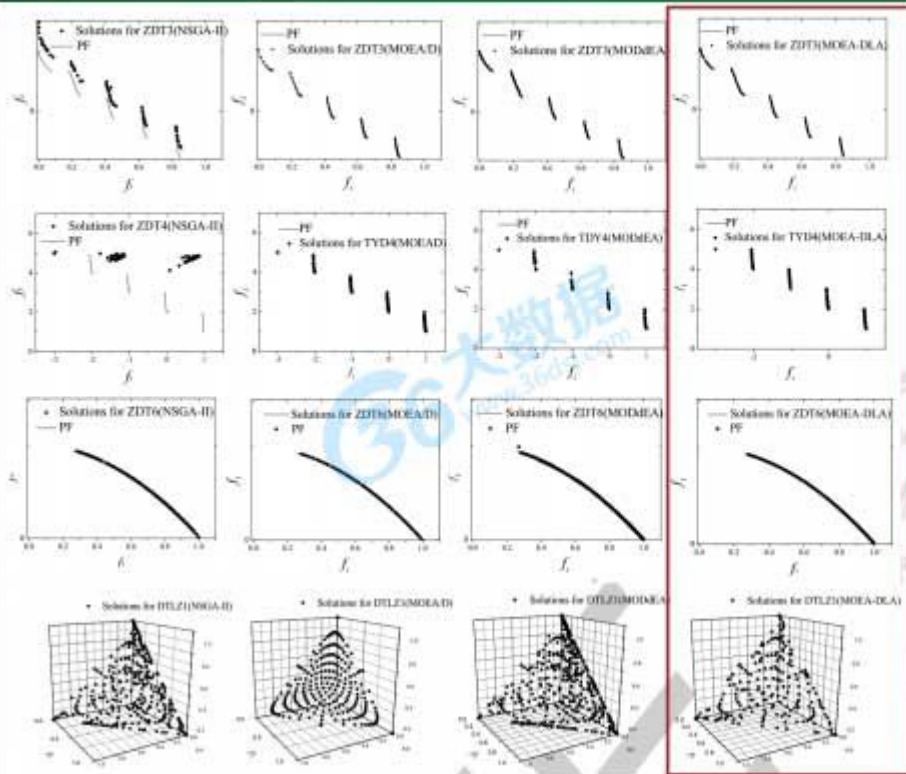




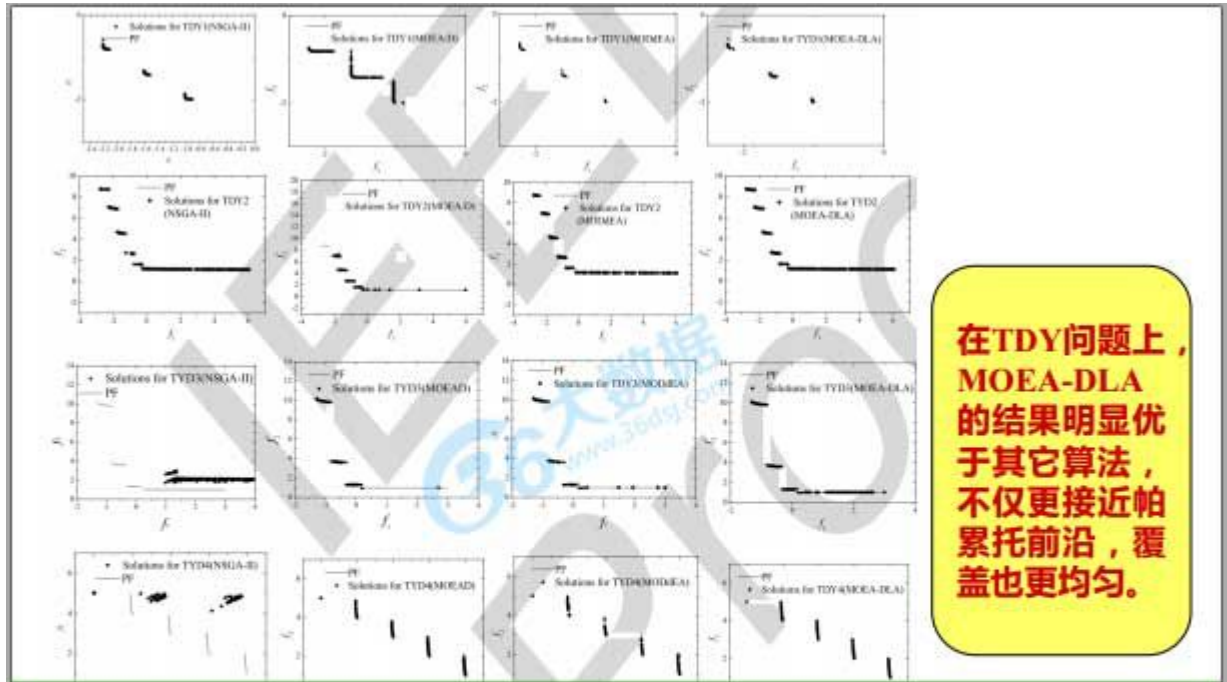
进化过程中个体的流向



研究2：带双层Archive的多目标进化算法



在ZDT, WFG
和DTLZ问题
上, 我们的算
法均能很好近
似帕累托前沿。



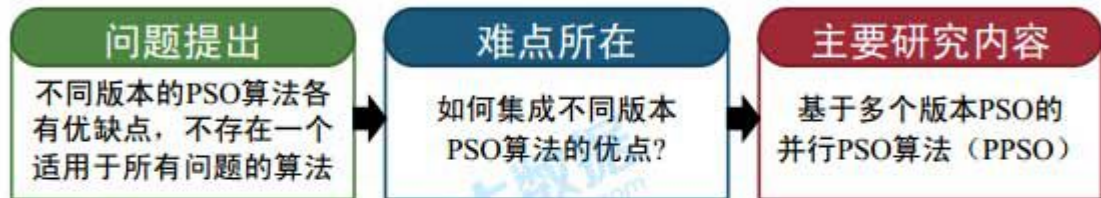
在TDY问题上，MOEA-DLA的结果明显优于其它算法，不仅更接近帕累托前沿，覆盖也更均匀。

研究成果已被IEEE Transactions on Cybernetics接收

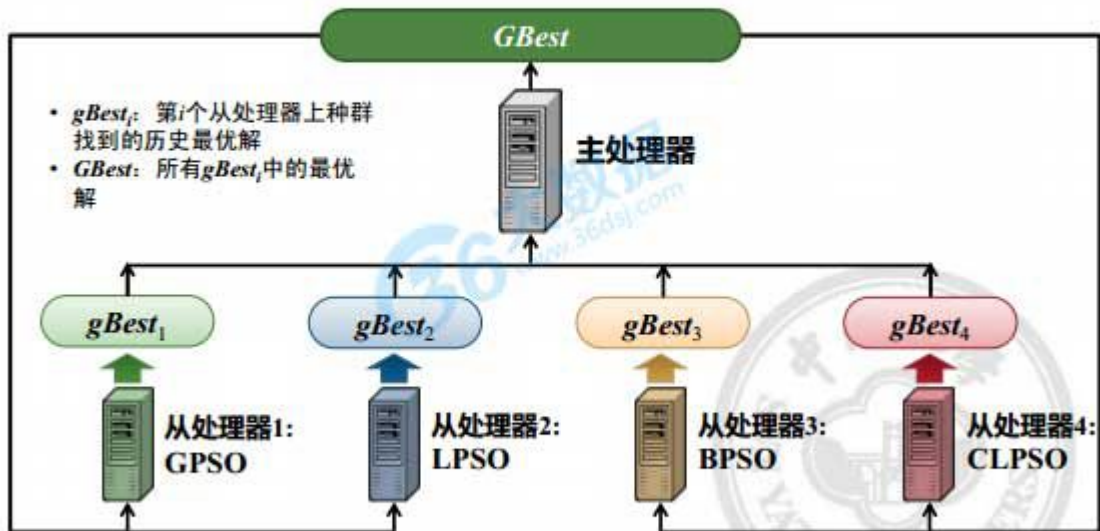
N. Chen, W.-N. Chen, Y.-J. Gong, Z.-H. Zhan, J. Zhang, Y. Li, & Y.-S. Tan, "An evolutionary algorithm with double-level archives for multiobjective optimization," in press.



研究背景与意义



算法思想



主要创新点：通过不同算法的并行执行和有效通信，提高算法的求解效率和普适性。

处理器间的通信机制

◎ 采用主-从拓扑结构：

从处理器 → **主处理器**

- 从处理器的进化过程被划分为若干阶段；
- 每阶段结束时，从处理器向主处理器发送其种群的历史最优解。

主处理器 → **从处理器**

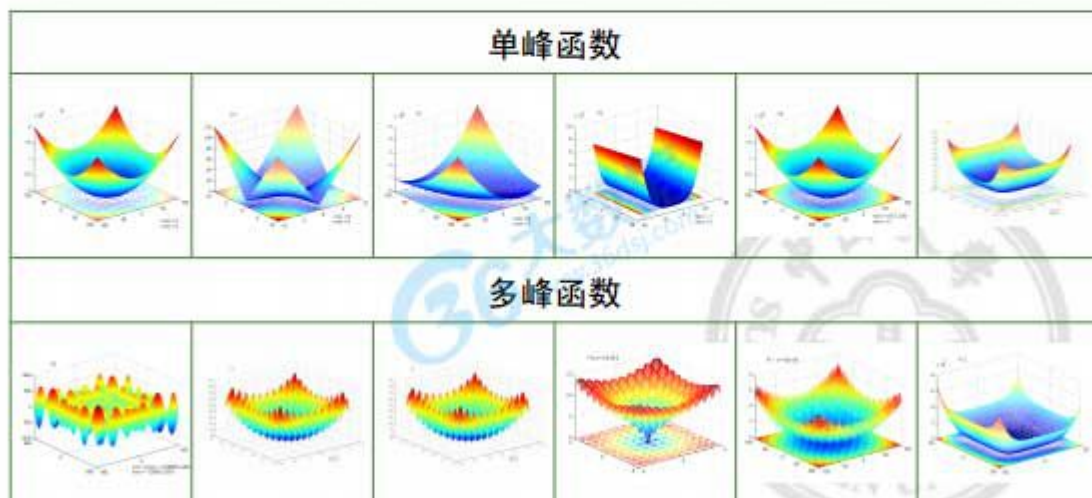
- 主处理器从所接收的解中挑选全局最优解，并将其发送给各个从处理器；
- 从处理器将其种群中随机选择的任意一个个体替换为全局最优解。

通讯机制的作用：

- ☑ 从处理器分阶段多次向主处理器传递最优解：**充分发挥不同算法的搜索能力**
- ☑ 主处理器将全局最优解返回给从处理器：**保持多样性，提供更好的引导向量**

实验结果与讨论

◎ 在12个30维标准测试函数上进行了测试



研究3：并行粒子群优化算法



实验结果与讨论

(30次运行的平均结果，**红色粗体**表示最佳结果，**蓝色斜体**表示次优结果)

| F | GPSO | LPSO | BPSO | CLPSO | PPSO |
|-------|------------------------|------------------------|---|------------------------|---|
| f_1 | 3.76×10^{-77} | 3.09×10^{-18} | 3.35×10^{-218} | 1.19×10^{-74} | 3.75×10^{-216} |
| f_2 | 2.00×10^{-41} | 2.36×10^{-13} | 2.39×10^{-152} | 4.16×10^{-47} | 9.33×10^{-152} |
| f_3 | 1.66×10^{-3} | 499.85 | 4.83×10^{-13} | 8.14×10^{-2} | 1.37×10^{-12} |
| f_4 | 0.22 | 9.39 | 7.41×10^{-05} | 0.27 | 1.48×10^{-05} |
| f_5 | 44.60 | 37.02 | 10.08 | 32.72 | 2.85 |
| f_6 | 0 | 0 | 3.1 | 0 | 0 |
| f_1 | 3089.33 | 4242.06 | 1342.33 | 43.44 | 1.34×10^{-2} |
| f_2 | 26.47 | 37.63 | 75.92 | 7.96 | 6.87 |

研究3的阶段性成果已发表于国际会议IES 2014

G.-W. Zhang, Z.-H. Zhan, K.-J. Du, Y. Lin, W.-N. Chen, J.-J. Li, & J. Zhang, "Parallel particle swarm optimization using message passing interface," in *Proceeding of IES 2014*, Singapore.

谢谢!



