

Алгоритм репрезентативного сэмплинга для параллельных систем баз данных*

Д.Д. Янцен, М.Л. Цымблер

Южно-Уральский государственный университет (Челябинск)

СЭМПЛИНГ

Сэмплинг базы данных представляет собой получение выборки из неё с целью оценить статистические параметры данных или ускорить вычисления в интеллектуальном анализе данных

Последовательные СУБД

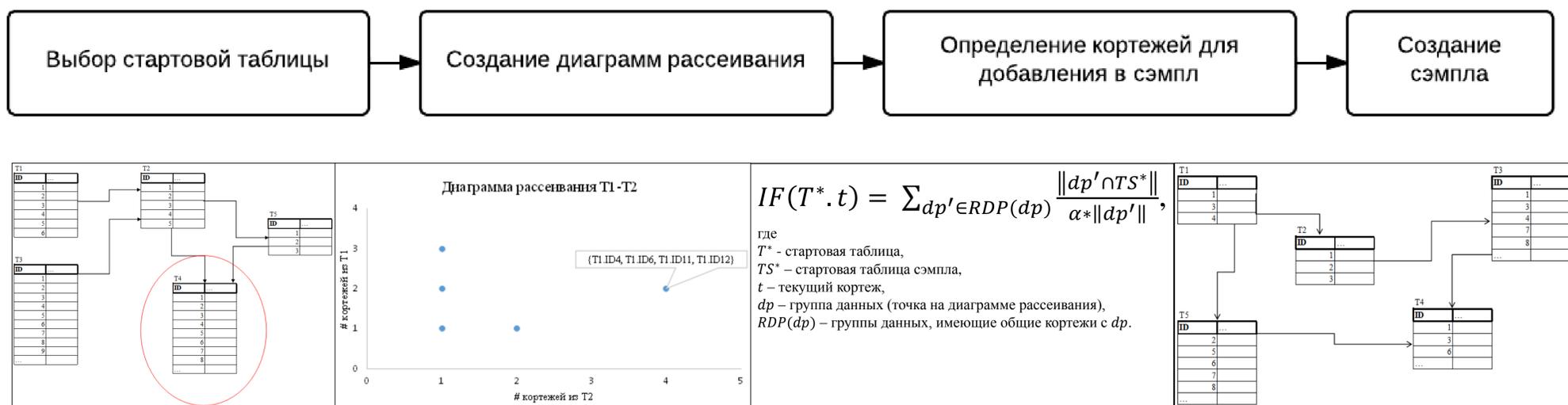
Случайный сэмплинг – выборка производится независимо от значений. В репрезентативном сэмплинге выборка сохраняет статистические особенности данных.

Параллельные СУБД

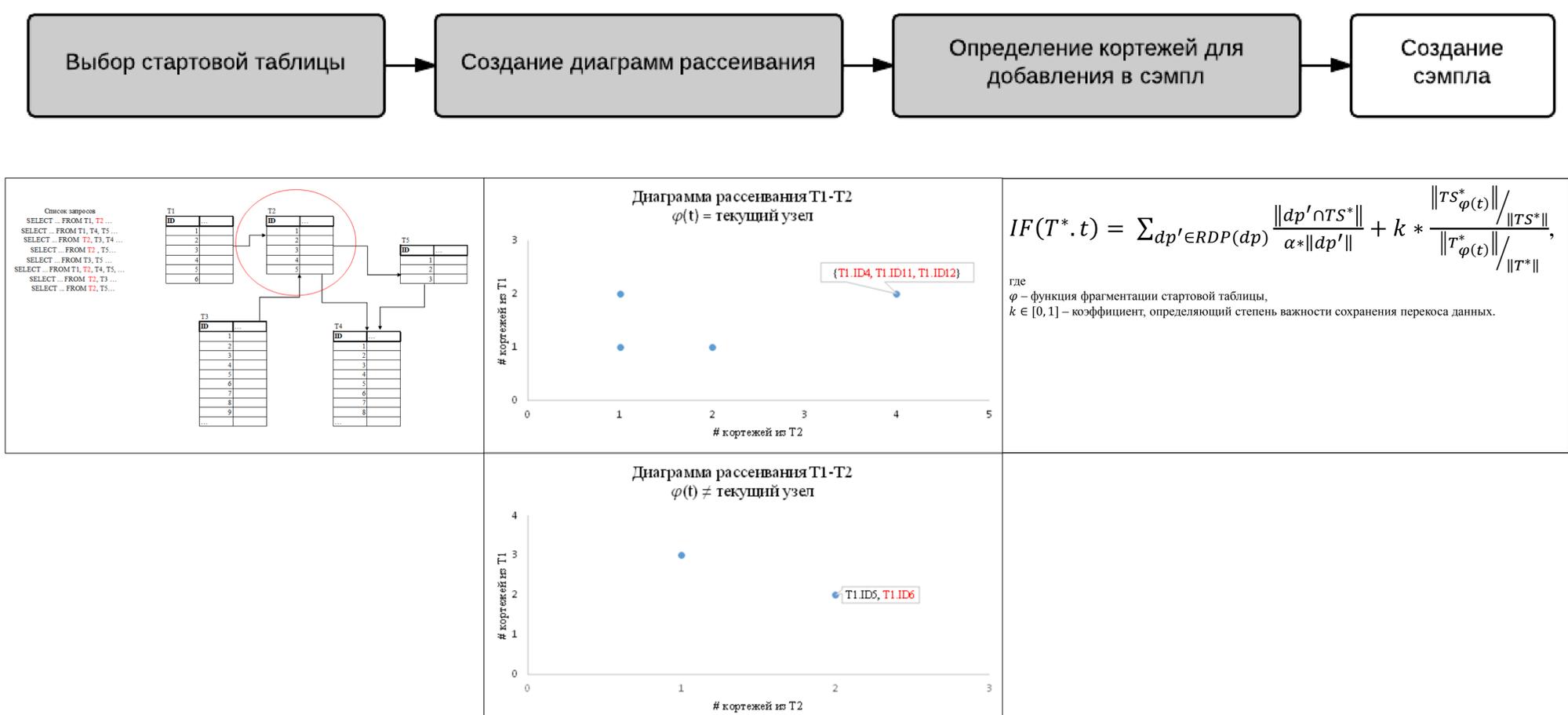
Основные проблемы выполнения сэмплинга в ПСУБД:

- Сохранение соотношения размеров фрагментов системы;
- Сохранение соотношения кортежей, которые необходимо передавать на другие узлы при выполнении запроса, к кортежам, которые должны быть обработаны на текущем узле.

Алгоритм CoDS для последовательных баз данных



Модификация алгоритма CoDS для ПСУБД



* Исследование выполнено при финансовой поддержке РФФИ в рамках научного проекта № 12-07-00443-а.