

# **ソフトバンク接続率 No.1 達成の裏で 活躍したビッグデータ解析と さらなる発展について**

**2015年6月3日**

**株式会社Agoop 代表取締役 柴山 和久、加藤 有祐**

**Redshift?**

# Why Redshift ?

シンプルSQL



初期費用不要



容量単価が安い



導入が早い



拡張性



操作性



目々 110 クエリ

解析 30 種類



4時間→2時間  
速度**2**倍

2.5T→16T  
容量**6.4**倍



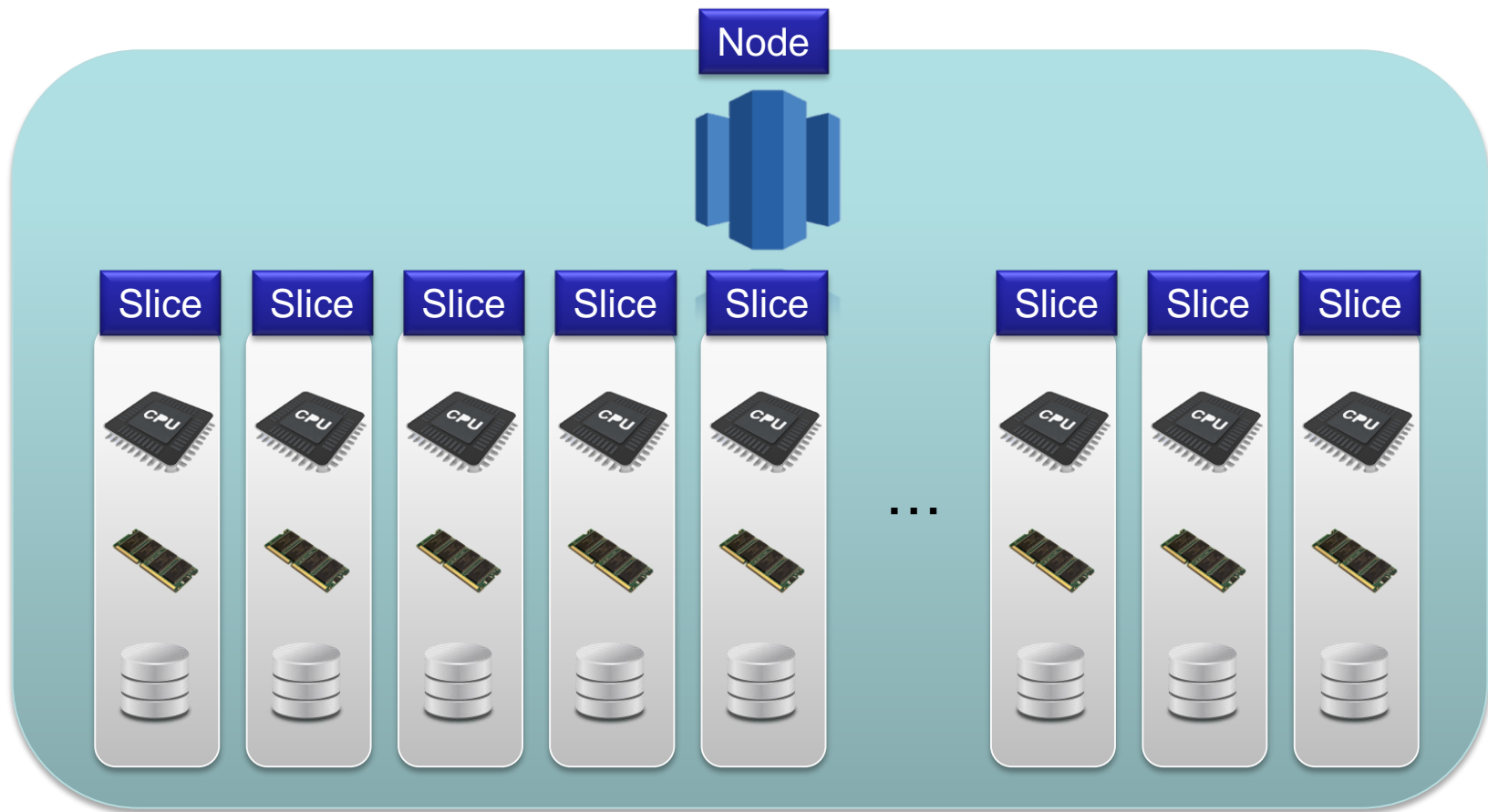
長期間のログ保存  
アドホック解析

短期間のログ保存  
日次解析



さらなる高速化へ

# 分散キーはパフォーマンスの最重要



## 分散キーの見直しで処理速度**2倍**！

user_id	action_id
1	2
2	2
2	3
2	5
5	2

**分散キー**を正しくセットしないと

データの**再移動**が発生

分散とクエリがマッチすると**ノード追加**の効果大

### 均等分散

1  
2

2  
5

2

### ALL分散

1  
2  
2  
2  
5

1  
2  
2  
2  
5

1  
2  
2  
2  
5

### キー分散

1

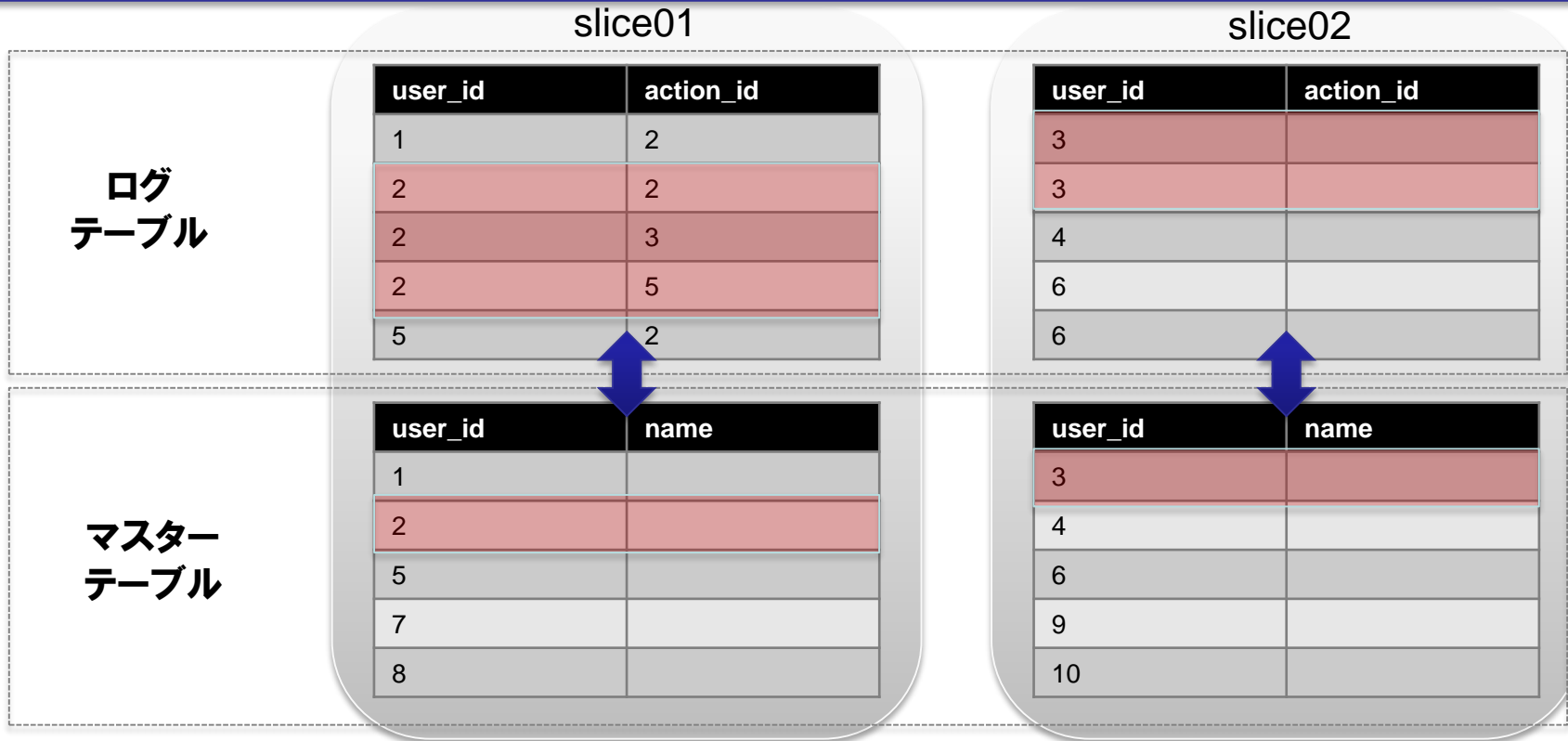
2  
2  
2

5

- ・ID系をセット
- ・JOINで使う項目
- ・Group Byする項目



## 結合キー, GroupByにセットする = 再分散しない



## DW1 x 2node クエリチューニング **2.5時間 → 40分**

### 大量データテーブルのIDを振り直す

Before

Column1(uniqueid)
Column2
Column3
....
Column39
Column40

**ノード間の大量データ移動発生**



**INSERT INTO  
SELECT**

row_number()
--------------

After

Column1(uniqueid)
row_number()



**Tempテーブル**

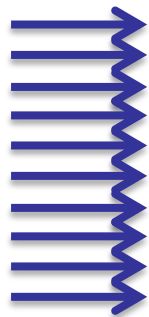
Column1(uniqueid)
Column2
Column3
....
Column39
Column40



**JOINしてINSERT**

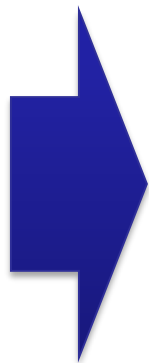


applicationlog01.gz  
applicationlog02.gz  
applicationlog03.gz  
applicationlog04.gz  
applicationlog05.gz  
applicationlog06.gz  
applicationlog07.gz  
applicationlog08.gz

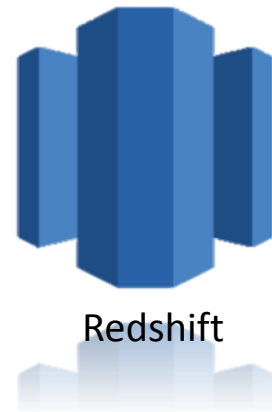


1ファイルずつCopy

applicationlog01.gz  
applicationlog02.gz  
applicationlog03.gz  
applicationlog04.gz  
applicationlog05.gz  
applicationlog06.gz  
applicationlog07.gz  
applicationlog08.gz



一括指定でCopy



**2000万件 (40GB) 1時間 → 15分へ短縮**