

---

# IMPROVING OPTIMALITY AND SPEED OF GREEDY GROUP RECURSION ALGORITHM

---

Florin Dobrian<sup>1</sup> Oleg Puzyrko White<sup>1</sup>

## ABSTRACT

Recently (Liu et al., 2025) has presented efficient algorithm - Greedy Group Recursion (GGR) - for reordering the rows and the fields within each row of an input table to maximize key-value (KV) cache reuse when performing LLM serving. In this paper, we propose several adjustments to GGR algorithm that can improve optimality of the solution and reduce its execution time.

## 1 INTRODUCTION

There has been growing research on LLM inference optimization. In particular, recent work (Liu et al., 2025; Cheng et al., 2025) presents solutions to optimize relational data analytics workloads for offline LLM inference. It proposes Greedy Group Recursion (GGR), an approximate algorithm that leverages functional dependencies (such as primary and foreign key relationships from the data schema) and table statistics, which are readily available in many databases and analytics systems, to reduce the search space.

The GGR algorithm is described in detail in (Liu et al., 2025). Let us recap briefly its main points.

The pseudocode specification of GGR is in Algorithm 1 and is taken from the original paper but with several critical corrections (and multiple little typo fixes).

## REFERENCES

Cheng, A., Liu, S., Pan, M., Li, Z., Agarwal, S., Cemri, M., Wang, B., Krentsel, A., Xia, T., Park, J., Yang, S., Chen, J., Agrawal, L., Naren, A., Li, S., Ma, R., Desai, A., Xing, J., Sen, K., Zaharia, M., and Stoica, I. Let the barbarians in: How ai can accelerate systems performance research, 2025. URL <https://arxiv.org/abs/2512.14806>.

Liu, S., Biswal, A., Kamsetty, A., Cheng, A., Schroeder, L. G., Patel, L., Cao, S., Mo, X., Stoica, I., Gonzalez, J. E., and Zaharia, M. Optimizing llm queries in relational data analytics workloads, 2025. URL <https://arxiv.org/abs/2403.05821>.

---

### Algorithm 1 Greedy Group Recursion (GGR)

---

```
1: Input: Table  $T$ , Functional Dependency  $FD$ 
2: Output: Prefix Hit Count  $S$ , Reordered List of Tuples  $L$ 

3: function HITCOUNT( $v, c, T, FD$ )
4:    $R_v \leftarrow \{i \mid T[i, c] = v\}$ 
5:   inferred_cols  $\leftarrow \{c' \mid (c, c') \in FD\}$ 
6:   inferred_vals  $\leftarrow \{T[R_v[1], c'] \mid c' \in \text{inferred\_cols}\}$ 
7:   tot_len  $\leftarrow \text{len}(v)^2 + \sum_{c' \in \text{inferred\_cols}} \left( \frac{\sum_{r \in R_v} \text{len}(T[r, c'])}{|R_v|} \right)^2$ 
8:    $HC \leftarrow \text{tot\_len} \times (|R_v| - 1)$ 
9:   cols  $\leftarrow [c] + \text{inferred\_cols}$ 
10:  vals  $\leftarrow [v] + \text{inferred\_vals}$ 
11:  return  $HC, cols, vals$ 
12: end function

13: function GGR( $T, FD$ )
14:   if  $|T|_{rows} = 1$  then
15:     return  $0, [T[1]]$ 
16:   end if
17:   if  $|T|_{cols} = 1$  then
18:      $S \leftarrow \sum_{v \in \text{distinct}(T[:, 1])} \text{HITCOUNT}(v, 1, T)$ 
19:     return  $S, \text{sort}([T[i] \mid i \in 1 \dots |T|_{rows}])$ 
20:   end if
21:    $b\_v, b\_c \leftarrow \text{None}, \text{None}$ 
22:    $max\_HC, b\_cols, b\_vals \leftarrow -1, [], []$ 
23:   for  $c \in \text{columns}(T), v \in \text{distinct}(T[:, c])$  do
24:      $HC, cols, vals \leftarrow \text{HITCOUNT}(v, c, T, FD)$ 
25:     if  $HC > max\_HC$  then
26:        $b\_v, b\_c \leftarrow v, c$ 
27:        $max\_HC, b\_cols, b\_vals \leftarrow HC, cols, vals$ 
28:     end if
29:   end for
30:    $R\_v \leftarrow \{i \mid T[i, b\_c] = b\_v\}$ 
31:    $A\_HC, A\_L \leftarrow \text{GGR}(T[\text{rows} \setminus R\_v, cols], FD)$ 
32:    $B\_HC, B\_L \leftarrow \text{GGR}(T[R\_v, cols \setminus b\_cols], FD)$ 
33:    $C\_HC \leftarrow max\_HC$ 
34:    $S \leftarrow A\_HC + B\_HC + C\_HC$ 
35:    $L \leftarrow [b\_vals + B\_L[i] \mid i \in 1 \dots |R\_v|] + A\_L$ 
36:   return  $S, L$ 
37: end function

38: return GGR( $T, FD$ )
```

---

<sup>1</sup>Data Analytics Group. Corresponding author: Oleg P. White <oleg.p.white@gmail.com>.