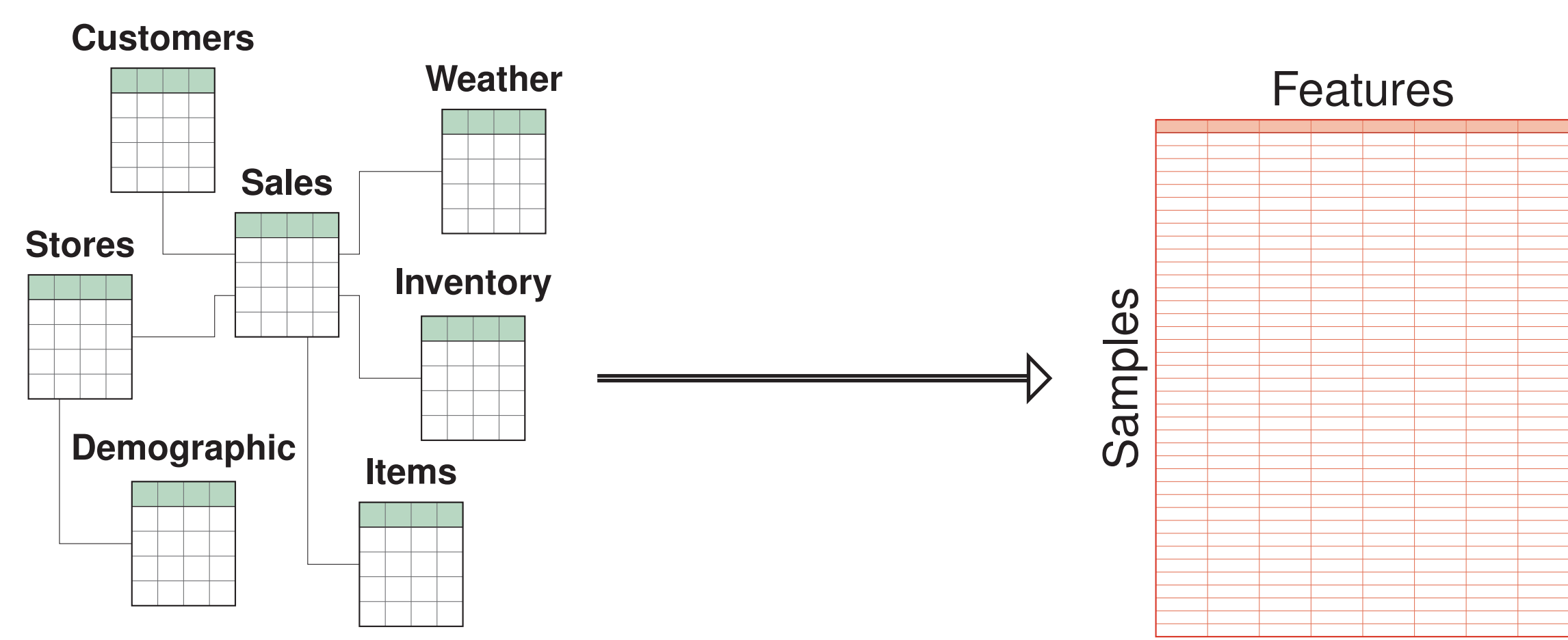


## Current State of Affairs

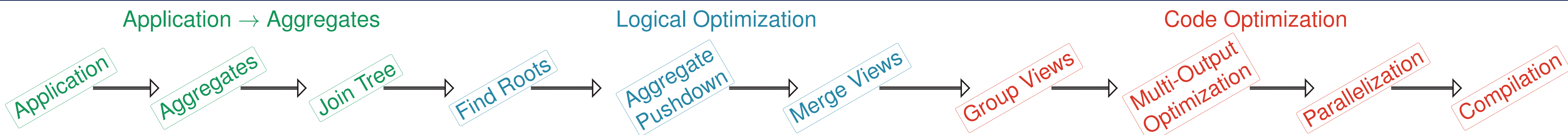


- Carefully crafted by domain experts
- Comes with relational structure
- Throws away relational structure
- Can be order-of-magnitude larger

## Turn Analytics Problem into Database Problem!

- Exploit structure in the data**
  - Algebraic structure: Factorized aggregate computation
  - Combinatorial structure: Query complexity measures
- Sharing computation and data access**
  - Aggregates decomposed into directional views over join tree
  - Share data access across views
- Specialization for workload and data**
  - Generate code specific to the query batch and dataset
  - Improve cache locality for hot data
- Parallelization**
  - Task and domain parallelism

## The Layers of LMFAO: Layered Multiple Functional Aggregate Optimization

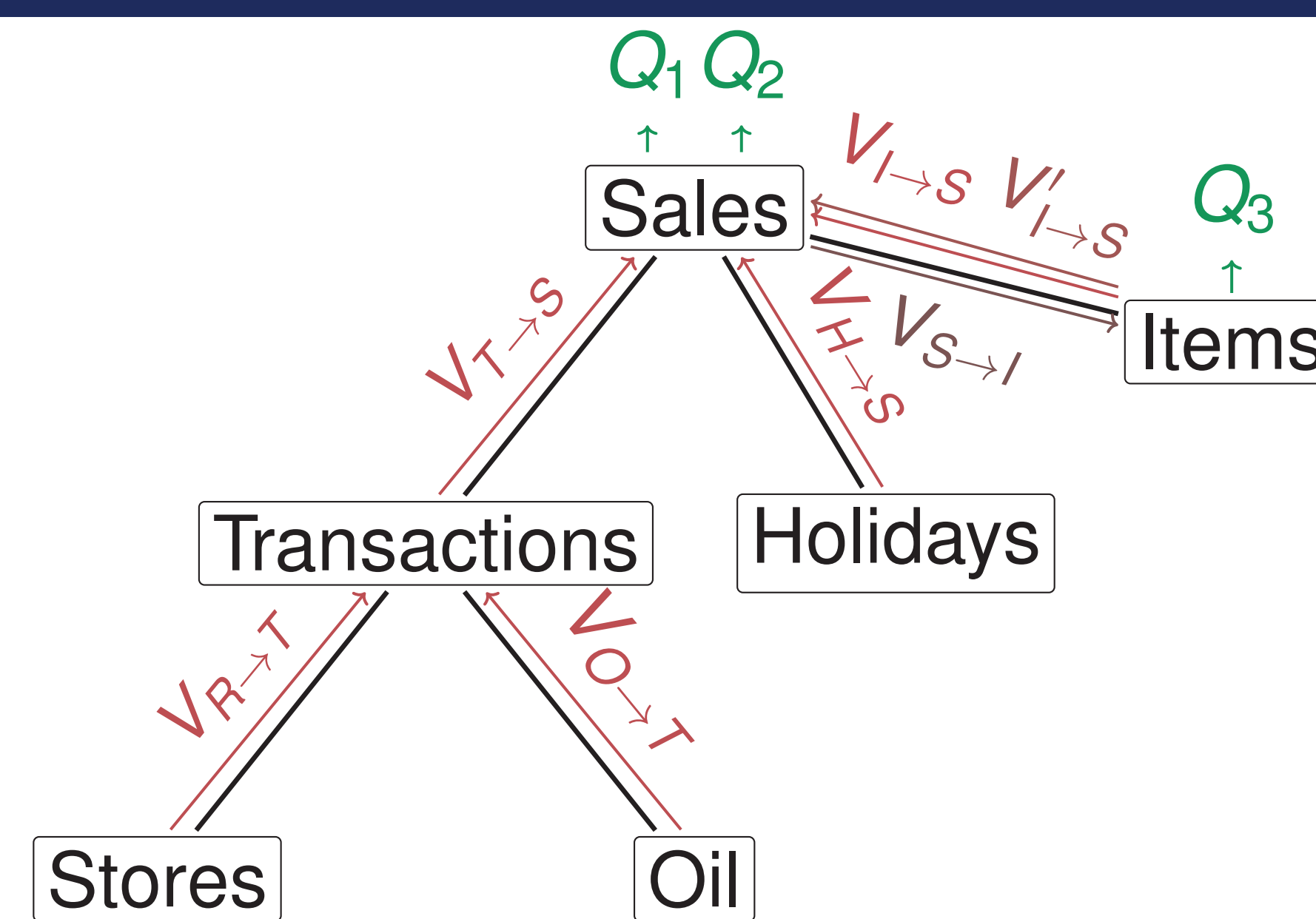


## Aggregates required per Analytics Workload

Workload	Query Batch	# Queries
Linear Regression	$SUM(X_i * X_j)$	814
Covariance Matrix	$SUM(X_i) GROUP BY X_j$ $COUNT(*) GROUP BY X_i, X_j$	
Decision Tree (Regression, 1 Node)	$VARIANCE(Y) WHERE X_j = c_j$	3,141
Mutual Information	$COUNT(*) GROUP BY X_i$	56
Chow-Liu Trees	$COUNT(*) GROUP BY X_i, X_j$	
Data Cubes	$SUM(M) GROUP BY X_1, \dots, X_d$	40

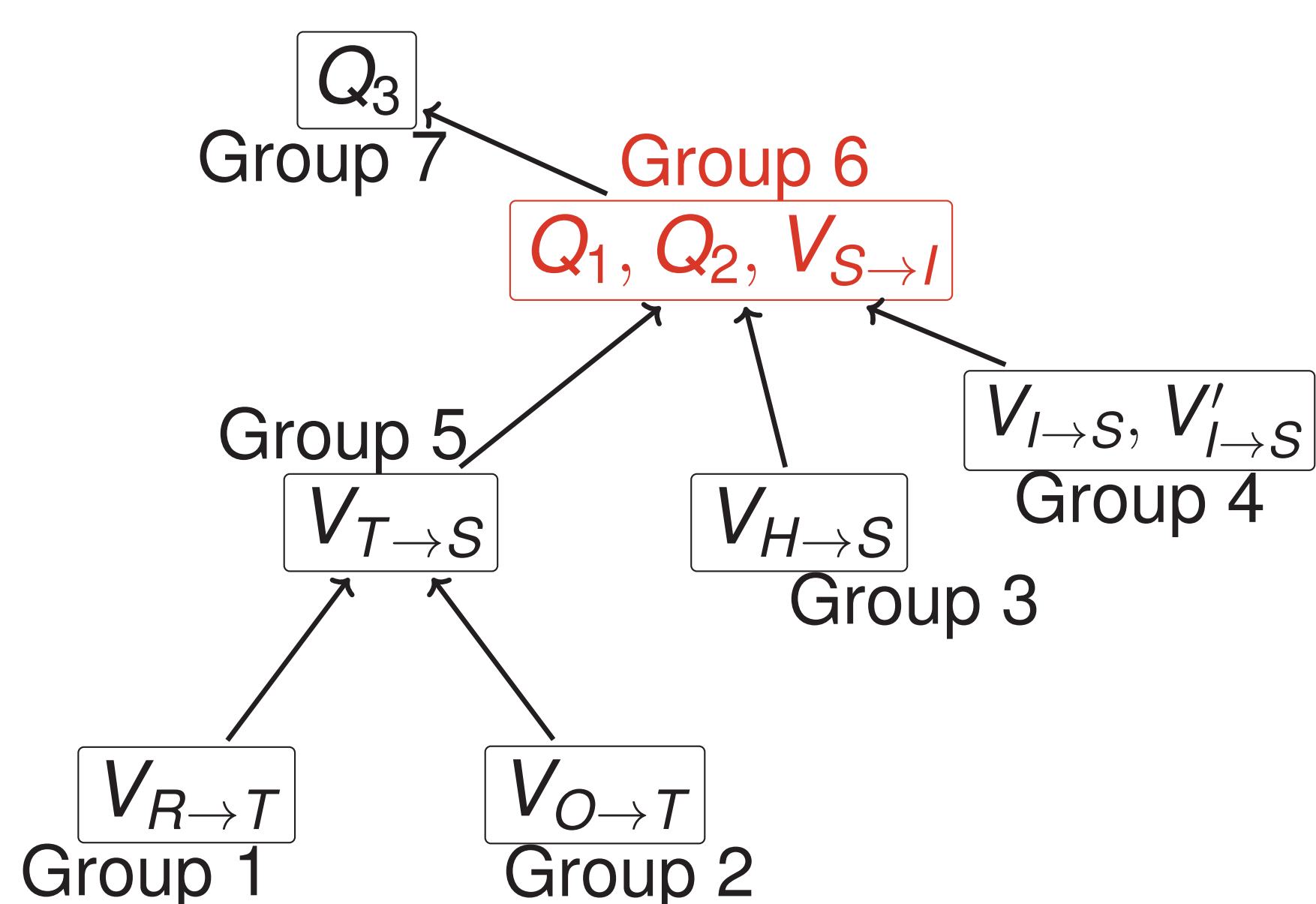
(# Queries shown for Retailer dataset)

## Logical Optimization



- Find Roots:** For each query, decide its output (root) node
- Break down each query into directional views over the join tree
- Reuse partial-aggregates and **Merge Views** with same group-by attributes

## Dependency Graph & View Groups



- Create Dependency Graph and **Group Views** computed over same relation
- View Group is computational unit, computed in one pass over relation
- Task and Domain Parallelism**

## Multi-Output Optimization & Code Compilation

```

V_I \to item
V'_I \to item
V_H \to date
V_T \to store
Q_1
Q_2
Q_3

```

```

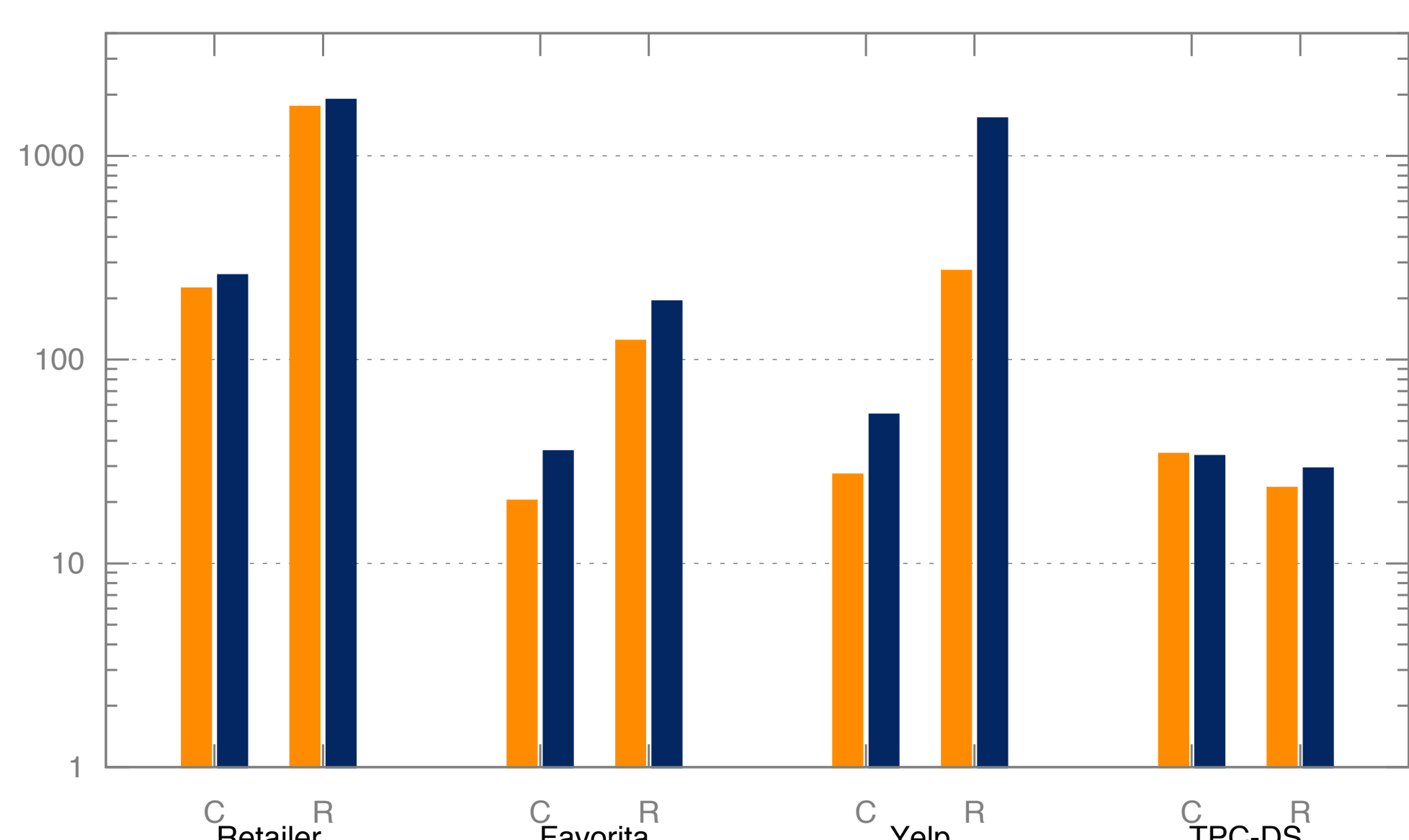
alpha_0 = 0;
foreach i in pi_item(S x_item V_I x_item V'_I)
  alpha_1 = V_I(i); alpha_2 = g(i); alpha_3 = 0;
  foreach d in pi_date(sigma_item=i S x_date V_H x_date V_T)
    alpha_4 = V_H(d); alpha_5 = 0;
    foreach c in pi_color(sigma_item=i V'_I : alpha_5 += h(d, c) * V'_I(i, c);
    alpha_6 = 0; alpha_7 = alpha_2 * alpha_5 * alpha_4;
    foreach s in pi_store(sigma_item=i, date=d S x_store sigma_date=d V_T)
      alpha_8 = V_T(d, s); alpha_9 = 0; alpha_10 = |sigma_item=i, date=d, store=s S|;
      foreach u in pi_units(sigma_item=i, date=d, store=s S : alpha_9 += f(u);
      alpha_6 += alpha_8 * alpha_9; alpha_11 = alpha_7 * alpha_8 * alpha_10;
      if Q_2(s) then Q_2(s) += alpha_11 else Q_2(s) = alpha_11;
    alpha_3 += alpha_4 * alpha_6;
    alpha_0 += alpha_1 * alpha_3; V_{S \to I}(i) = alpha_3 * alpha_2;
Q_1 = alpha_0;

```

$Q_1$ :  $SUM(f(units))$        $V_{S \to I}$ :  $SUM(f(units) \cdot g(item))$  GROUP BY item  
 $Q_2$ :  $SUM(g(item) \cdot h(date, color))$  GROUP BY store

## Aggregate Experiments

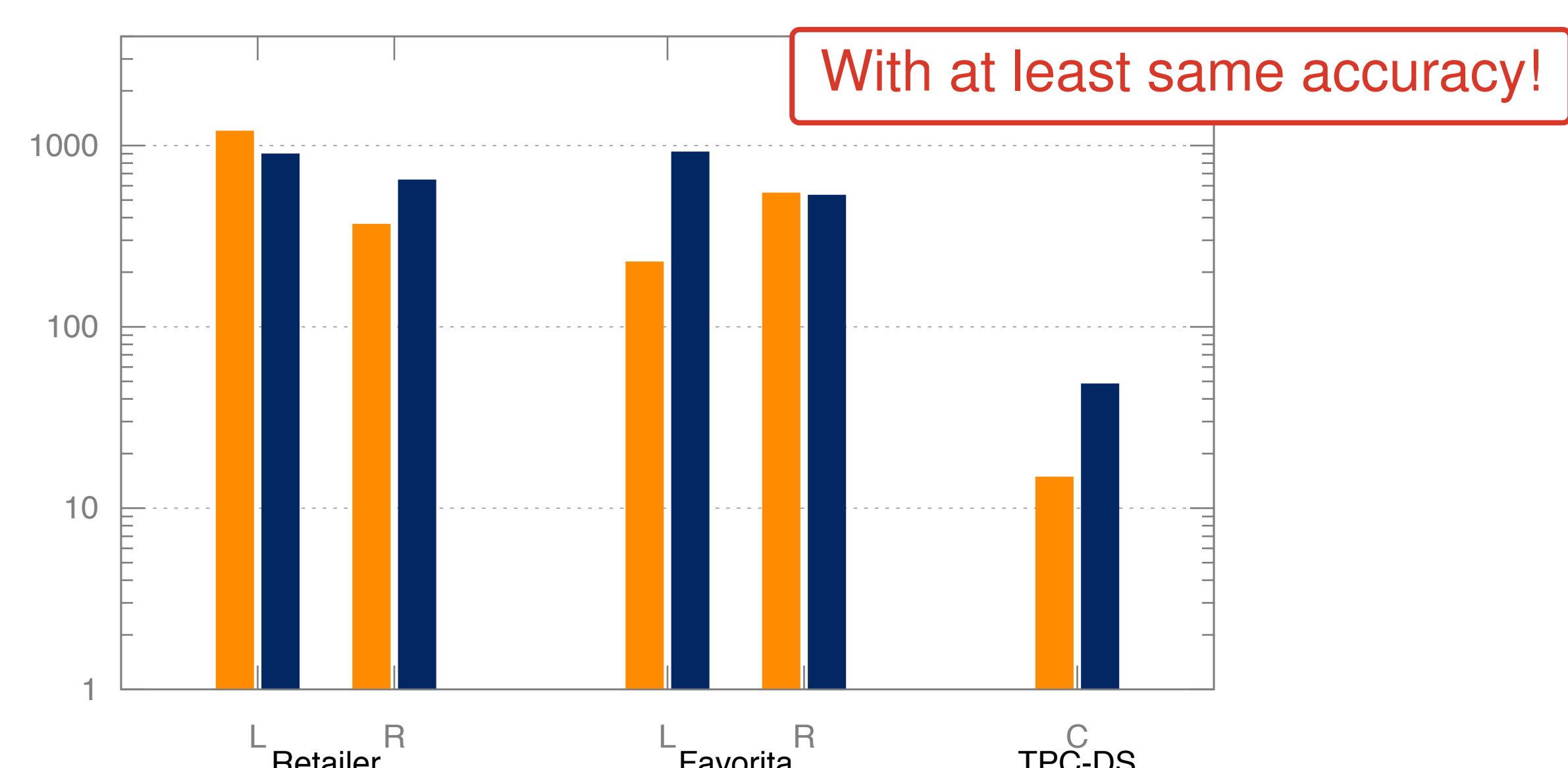
Relative Speedup over **DBX** and **MonetDB**



C = Covariance Matrix; R = Regression Tree Node; AWS d2.xlarge (4 vCPUs, 32GB)

## Machine Learning Experiments

Relative Speedup over **TensorFlow** and **MADlib**



L = Linear Regression; R = Regression Tree; C = Classification Tree; Intel i7-4770 (8 CPUs, 32GB)